

## Control Generation by Program Transformation

**Andy King and Jonathan C. Martin\***

*University of Kent, UK*

*a.m.king@kent.ac.uk*

---

**Abstract.** The objective of control generation in logic programming is to derive a computation rule for a program that is efficient and yet does not compromise program correctness. Progress in solving this fundamental problem in logic programming has been slow and, to date, only partial solutions have been proposed. Previously proposed schemes are either inefficient, incomplete (incorrect) or difficult to apply for programs consisting of many components (the scheme is not modular). This paper shows how the control generation problem can be tackled by program transformation. The transformation relies on information about the depths of derivations to derive delay declarations which orchestrate the control. To prove correctness of the transformation, the notion of semi-delay recurrency is introduced, which generalises previous ideas in the termination literature for reasoning about logic programs with delay declarations. In contrast to previous work, semi-delay recurrency does not require an atom to be completely resolved before another is selected for reduction. This enhancement permits the transformation to introduce control which is flexible and relatively efficient.

**Keywords:** logic programming, program transformation, control generation

### 1. Introduction

A logic program can be considered as consisting of a logic component and a control component [22]. Although the meaning of the program is largely defined by its logical specification, choosing the right control mechanism is crucial in obtaining a correct and efficient program. One of the most popular ways of defining control is via suspension mechanisms which delay the selection of an atom in a goal until some condition is satisfied [9, 32]. Delays have been applied to, among other things, handle negation [13], delay non-linear constraints [21], support coroutining [40], improve search by rescheduling tests [12], implement concurrency protocols [39] and induce termination [33]. In general, delay mechanisms

---

Address for correspondence: Andy King, Computing Laboratory, University of Kent, Canterbury, CT2 7NF, UK

\*The authors gratefully acknowledge Nuffield grant SCI/180/94/G which funded their collaboration

are used to define dynamic selection rules with the two main aims of enhancing performance through coroutining or expressiveness whilst ensuring termination. However, reasoning about logic programs with delay is notoriously difficult and often completeness, and hence program correctness, is inadvertently sacrificed by the programmer in the quest for improved efficiency or generality.

Consider, for example, the Prolog program that is listed below:

```
inorder(nil, []).
inorder(tree(L,V,R), I) :- inorder(L, LI), inorder(R, RI), append(LI, [V|RI], I).

append([], Xs, Xs).
append([X|Xs], Ys, [X|Zs]) :- append(Xs, Ys, Zs).
```

Declaratively, the program defines the relation that the second argument (a list) is an in-order traversal of the first argument (a tree). With leftmost selection – the default computational rule of Prolog – the goal `inorder(tree(nil, a, tree(nil, b, nil)), L)` will terminate and bind `L` to `[a, b]`. The goal `inorder(T, [a, b])` will also return an answer, namely, `T = tree(nil, a, tree(nil, b, nil))`. One subtlety, however, is that this goal has another answer, namely, `T = tree(tree(nil, a, nil), b, nil)` yet a request to compute this answer results in a diverging computation. Thus `inorder` existentially terminates (it finds one answer) but it does not universally terminate (it does not enumerate all the answers) [14]. The program can be enhanced by adding the delay declaration `:- block append(-, ?, -)` which asserts that an `append` atom cannot be selected until either its first or third argument is instantiated to a non-variable term. (Note that since the syntax and semantics of delay declarations have not been standardised [16], for purposes of exposition the paper adopts the block declarations of SICStus [37] since these can be simulated with, say, the `DELAY` declarations of Gödel [19] and the `wait` declarations of MU-Prolog [32].) This block declaration, however, is not by itself sufficient for universal termination in both modes. In addition, the body atoms of `inorder` have to be reordered to obtain:

```
inorder(nil, []).
inorder(tree(L,V,R), I) :- append(LI, [V|RI], I), inorder(L, LI), inorder(R, RI).

:- block append(-, ?, -).
append([], Xs, Xs).
append([X|Xs], Ys, [X|Zs]) :- append(Xs, Ys, Zs).
```

Then both goals `inorder(tree(nil, a, tree(nil, b, nil)), L)` and `inorder(T, [a, b])` universally terminate (at least with a computation rule that invokes a suspended goal as soon as its delay condition is satisfied [37]). This program illustrates an unsatisfactory feature of logic programming: although the declarative behaviour of `inorder` is transparent, its operational behaviour is opaque because of the control. The program also illustrates that delay declarations *can* be applied to obtain completeness and correctness. The problem is how to add delays systematically; the challenge is how to add delays automatically.

The seminal paper on logic and control [22] has been interpreted by many as setting an agenda for logic programming: control generation [8, 20, 24, 33, 34, 42]. The objective of control generation in logic programming is to derive a computation rule for a program that is efficient and yet does not compromise completeness, hence universal termination. Progress in solving this fundamental problem has been slow, despite recent advances on termination checking in the context of sophisticated control [6, 35]. However,

crucial insight into this problem is given by Naish [32] within a discussion on the desirable properties of computation rules. Naish writes,

“There is a slightly more subtle rule [computation rule tactic] which applies more to goals which have solutions ... Although the success branches of a tree are fixed, the number and length of the other branches are not. The rule, therefore, is to avoid the creation of unnecessary failure (and infinite) branches”.

Martin and King [28] introduce a program transformation which applies this tactic to generate a control component from the logic component that prunes out infinite branches. Specifically, if the maximum depth of the SLD-tree needed to solve a given goal can be determined, then by only searching to that depth, the goal will be completely solved. Hence all the answers will be obtained, in a finite number of steps, thereby enforcing termination without compromising completeness.

Pedreschi and Ruggieri [34] apply the same tactic with a similar transformation. This latter transformation adds a counter that is decremented on each resolution step. Once the counter reaches zero, the derivation is failed. The success of the transformation is predicated on being able to find an upper bound on the depth of all the successful, finite SLD-trees. By way of contrast, the transformation advocated by Martin and King [28] introduces not one but possibly many counters; one counter is used for each recursive loop in the program to control how many times it is traversed. One advantage of this approach is that it gives finer control on the depth to which the search space is probed, thereby improving efficiency. Moreover, it is easier to compute an upper bound on how many times a loop is traversed than find an upper bound on the total depth of the SLD-tree. The issue here is one of modularity [2]. Ideally a termination argument or dually, a construction that enforces termination, should be modular in the sense that it should be derived in a bottom-up fashion by considering, in turn, each of the strongly connected components (SCCs) in the predicate dependency graph. The problem that complicates the derivation of the level mappings [29] which underpin termination analysis is the same problem that complicates the derivation of a depth bound in control generation. In both cases, it is necessary to consider not only the recursive behaviour of those predicates in the top-level SCC, but those predicates (sub-computations) that are invoked from within it. The modularity problem in termination analysis is to synthesise a level mapping which ensures that each call to a sub-computation is bounded. The modularity problem in control generation is to infer a depth bound such that each sub-computation falls within the bound. The approach to control generation proposed in this paper alleviates much of this problem.

Since this paper draws together a number of threads in the termination and control generation literature, as well as extending preliminary work on this topic [28], we summarise its contributions:

- The paper explains how the control generation problem can be tackled by program transformation. It also shows that the performance of a program developed with control generation can compare favourably with that of a logically equivalent program equipped with *ad hoc* delay declarations.
- The paper explains how control generation can be tackled in a modular fashion.
- The paper formally argues the correctness of the transformation. Since it introduces delays to orchestrate the control, the correctness proof is non-trivial. In fact, the argument is formulated in terms of a new class of logic program – the class of semi-delay recurrent program (which was originally sketched in [28]) – which refines the concept of delay recurrency [26, 27]. Delay recurrency was proposed for proving termination of logic programs with delay. Delay recurrency is

sufficient for termination provided the program is executed under local selection whilst satisfying the delay conditions. However, under local selection, the selected atom is completely resolved, that is, those atoms it directly and indirectly introduces are also resolved, before any other atom is selected [43]. As Marchiori and Teusink [27] concede, this bars coroutining. By way of contrast, semi-delay recurrent programs permit coroutining and thereby provide a foundation for more flexible and more efficient control generation. Moreover, semi-delay recurrent programs can be directly implemented in systems such as SICStus Prolog [37] which implement the scheduling tactic that suspended goals should be invoked as soon as their delay conditions are satisfied [12].

- Finally, to make the ideas more accessible to a general program transformation audience, the paper reviews the problems that arise in control generation. It surveys some of the proposed solutions as well as their short-comings.

The review constitutes section 2. Section 3 illustrates the transformation tactic with a familiar example program. Section 4 summaries the key ideas in delay recurrency [27], before moving on to refine these concepts to introduce the class of semi-delay recurrent logic. Section 5 presents a formal development of the proposed transformation, including the correctness results. The transformed programs are by construction semi-delay recurrent and hence termination is guaranteed. Section 6 discusses automation and appraises of the performance and practicality of this transformational approach to control generation. Finally, section 7 presents the concluding discussion.

## 2. Issues of dynamism in control generation

The presence of delayed goals in a computation significantly complicates the termination behaviour of a logic program. This section reviews the subtleties which can arise when delays are used to induce termination, explaining why there is room for improving the solutions that have been previously proposed.

### 2.1. The issue of local boundedness

Consider the append program introduced previously, complete with a block declaration which delays the selection of an append atom until either the first or third arguments are non-variable:

```
:- block append(-, ?, -).
append([], Xs, Xs).
append([X|Xs], Ys, [X|Zs]) :- append(Xs, Ys, Zs).
```

Although the block declaration is intended to assist termination, it is not sufficient to ensure that *all* append goals terminate. As Naish [33] points out, the goal `append([X|Xs], Ys, Xs)` satisfies the condition in the declaration and therefore does not suspend, yet its derivation is an infinite one. The derivation is infinite since unifying the goal with the head of the second clause will instantiate `Xs` to `[X|Zs]` so the sub-goal `append(Xs, Ys, Zs)` is then an instance of the initial goal.

Termination can only be guaranteed for all goals by strengthening the condition in the delay declaration. This is where the trade-off between efficiency, termination and non-suspension takes place. Ideally, the delay declarations would be as weak as possible since strengthening a delay declaration amounts to inspecting more sub-terms of a goal which in turn decreases efficiency. On the other hand, the delay

condition should be strong enough to bar the resolution of goals which have infinite derivations. Yet, if the delay condition is too strong, goals may suspend which actually possess finite derivations. Thus one of the main problems in generating control is that of finding a suitable delay condition for a predicate that suspends those goals that have an infinite derivation, yet does not suspend goals which possess finite derivations. Moreover, the delay condition for a predicate should ideally be inexpensive to check. This trade-off primarily relates atomic goals – a call to a predicate and those it invokes – rather than the interaction between the sub-goals of a goal which shape how a search tree is traversed. The trade-off thus relates to *local* effects rather than global interactions. In particular, the trade-off relates to the degree of instantiation that arguments of an atomic goal require, or adopting terminology more in tune with the termination literature [14], whether the arguments of the goal are *bounded* with respect to some level mapping. This simply means that every instance of the goal is of fixed size where size is defined by a level mapping that maps goals to natural numbers. Thus, henceforth, this trade-off in control generation will be referred to as the *local boundedness* issue. There have been several attempts at addressing this problem, each of which will be discussed below with reference to the append program.

### 2.1.1. On using linearity

Lüttringhaus-Kappel [24] observed that for atomic goals, linearity is a property that is often sufficient for termination. Linearity relates to the number of times variables occur within a goal: a goal is non-linear if at least one variable occurs multiply within the goal; otherwise it is linear. Thus the goal  $\text{append}([X|Xs], Ys, Xs)$  is non-linear because the first and third arguments both contain the variable  $Xs$ . Delaying the goal  $\text{append}([X|Xs], Ys, Xs)$  until it is linear is equivalent to delaying the goal until  $Xs$  is instantiated to a term that contains no variables. This would indeed prevent looping. However, checking for linearity requires all the sub-terms within the arguments of a goal to be completely traversed and is therefore potentially expensive. Moreover, there is little prospect of applying abstract interpretation techniques to detecting linearity at compile-time since these analysis techniques rely on the control being predetermined. More significantly still, delaying a goal until it is linear would suspend goals with finite derivations such as  $\text{append}([X, X], Ys, Zs)$ . Thus linearity would be potentially expensive to enforce and would additionally suspend goals that have finite derivations.

### 2.1.2. On using rigidity

Mesnard [30] proposes generating control by delaying goals until their arguments are bound to rigid terms. A term is rigid with respect to a norm – a function that assigns a size to the term – if its size cannot change even when the variables within the term are further instantiated. For example, a term is rigid with respect to the list-length norm if the term is a list of determinate length. Mesnard [30] infers that  $\text{append}$  will terminate if called when its first or third arguments are rigid with respect to the list-length norm. Hence, calls to  $\text{append}$  are blocked until this rigidity property holds. This promising tactic successfully enforces termination but at the price of traversing (at least) one list on *each* call to the predicate. In fact, the check is only performed on the initial call [30], but no justification for this refinement is given. Thus the issues with this approach are primarily those of efficiency and correctness.

### 2.1.3. On using modes

Naish [33] proposes solving the problem of the goal `append([X|Xs], Ys, Xs)` with the use of modes which flag the fact that output should not feed back into the input. The mode system that Naish advocates is unusual in that modes are usually specifications on single procedures rather than conjunctions (a notable exception being the mode system of IC-Prolog which applied a form of mode to specify eager consumers and lazy producers [12]). However, although modes may form a good basis for solving this problem, they have not been shown to be satisfactory for reasoning about another termination problem, that of speculative output bindings. This is discussed below.

## 2.2. The issue of global boundedness

Even when finite derivations exist, delay conditions alone are not, in general, sufficient to ensure termination. Infinite computations may arise as a result of so-called speculative output bindings [33], which can arise from *global* interactions between the sub-goals of a goal due to dynamic selection. The effect that they have on termination is the focus of interest here and henceforth this issue will be referred to as the *global boundedness* issue. To illustrate the problem caused by speculative output bindings consider the `qsort` program listed in figure 1. The delay declarations on `leq` and `gr` postpone the tests  $X \leq Y$  and  $X > Y$  so as to avoid instantiation errors. The declarations on `qsort`, `part` and `append`, on the other hand, are an attempt to enrich the program so that it can be used in reverse mode, for example, with the call `qsort(X, [1, 2, 3])`. The declarations for `qsort`, `part` and `append` serve to illustrate the issue of global boundedness. The issue is illustrated by augmenting the program with the clause:

```
append(Xs, [_|Xs], Xs) :- fail.
```

The declarative semantics of the program are completely unchanged by the addition of this clause and one would hope that the new program would produce exactly the same set of answers as the original. But this is not the case as the goal `qsort(X, [1, 2, 3])` illustrates. Following resolution with the second clause of `qsort`, the only goal that can be selected is `append(Ls, [X|Bs], [1, 2, 3])`. When this is unified with the above clause, both `Ls` and `Bs` are bound to the list `[1, 2, 3]`. These bindings are speculative because if the sub-goal `fail` were to be selected, then they would be retracted. As a result of these speculative output bindings the previously suspended sub-goals `qsort(L, Ls)` and `qsort(B, Bs)` can resume *before* the `fail` is selected. The net result is an infinite derivation due to recurring sub-goals of the form `qsort(X, [1, 2, 3])`. The problem is that bindings are made before failure is detected and no matter how stringent the delay conditions, loops of this kind cannot generally be avoided without regard for the computation as a whole.

### 2.2.1. On using local selection rules

To remedy this problem, Marchiori and Teusink [26, 27] propose the use of a local selection rule. Such a rule only selects atoms from those that are most recently introduced in a derivation [43]. This ensures that any atom selected from a goal is completely resolved before any other atom in the goal is selected. The effect in the above example is that the `append` sub-goals would be fully resolved, hence `fail` would be selected before the `qsort` sub-goals are reawakened. This would prevent an infinite loop, but at the expense of not allowing any form of coroutining. This is a severe restriction since coroutining is a useful programming tactic within logic programming.

```

:- block qsort(-, -).
qsort([], []).
qsort([X|Xs], Ys) :- part(Xs, X, L, B), qsort(L, Ls), qsort(B, Bs), append(Ls, [X|Bs], Ys).
:- block part(-, ?, -, ?), part(-, ?, ?, -).
part([], -, [], []).
part([X|Xs], Y, [X|Ls], Bs) :- leq(X, Y), part(Xs, Y, Ls, Bs).
part([X|Xs], Y, Ls, [X|Bs]) :- gt(X, Y), part(Xs, Y, Ls, Bs).
:- block append(-, ?, -).
append([], Xs, Xs).
append([X|Xs], Ys, [X|Zs]) :- append(Xs, Ys, Zs).
:- block leq(-, ?), leq(?, -).
leq(X, Y) :- X ≤ Y.
:- block gt(-, ?), gt(?, -).
gt(X, Y) :- X > Y.

```

Figure 1. The quicksort algorithm expressed in Prolog

Local selection is also applied in the control generation work of Hoarau and Mesnard [20]. This work builds on the termination inference methods developed by Mesnard and his colleagues [31] and infers initial modes for a query that, if satisfied, ensure that a logic program left-terminates. The chief advance described in [20] is that the paper additionally infers how goals can be statically reordered so as to obtain termination with leftmost selection. Of course, leftmost selection, prohibits coroutining.

### 2.2.2. On delaying output unification

Naish [33] proposes finessing the problem of speculative output bindings by postponing output unification. In the case of the above example, the clause could be rewritten to:

```
append(Xs, Ys, Zs) :- fail, Ys = [_|Xs], Zs = Xs.
```

This transformation presupposes left-to-right scheduling as the default. The intended effect of the transformation is that no output bindings should be made until the computation is known to succeed. This tactic restricts coroutining since the act of postponing an output binding will inevitably postpone the activation of any goal waiting on that binding.

## 3. Worked example

This paper shows how program transformation can be applied to solve both the local and the global boundedness problem. The transformation introduces delay declarations that realise rigidity checks. These solve the local boundedness problem in the first instance by ensuring that the initial goal is bounded; sufficiently instantiated to assure termination. Boundedness of subsequent goals is ensured by depth counters that are introduced by the transformation so as to eliminate subsequent rigidity checks.

```

:- block msort(-, -).
msort([], []).
msort([X], [X]).
msort([X, Y|Xs], S) :- split(Xs, L1, L2), msort([X|L1], S1), msort([Y|L2], S2), merge(S1, S2, S).
:- block split(-, -, -).
split([], [], []).
split([X|L], [X|L1], L2) :- split(L, L2, L1).
:- block merge(-, ?, -), merge(?, -, -).
merge([], Ys, Ys).
merge(Xs, [], Xs).
merge([X|Xs], [Y|Ys], [X|Zs]) :- leq(X, Y), merge(Xs, [Y|Ys], Zs).
merge([X|Xs], [Y|Ys], [Y|Zs]) :- gt(X, Y), merge([X|Xs], Ys, Zs).

```

Figure 2. The mergesort algorithm expressed in Prolog

The global boundedness problem is also neatly solved. The depth counters ensure that the search space is finite, so even though speculative output bindings may still occur, they cannot lead to infinite derivations. Moreover, the global boundedness problem is overcome without grossly impeding coroutining.

To illustrate the approach, consider the mergesort program listed in figure 2, assuming the predicates `append`, `leq` and `gt` are defined as in figure 1. The delay declarations are an attempt to enhance the program so that it can be additionally used in reverse; the query `msort(L, [1, 2, 3])` should in principle terminate and systematically enumerate the solutions  $L = [1, 2, 3]$ ,  $L = [1, 3, 2]$ ,  $L = [2, 1, 3]$ ,  $L = [2, 3, 1]$ ,  $L = [3, 1, 2]$  and  $L = [3, 2, 1]$  though not necessarily in that order. The delay declaration `:- block merge(-, ?, -), merge(?, -, -)` asserts that a call to `merge` should block until (1) its first *or* third arguments are instantiated *and* (2) its second *or* third arguments are instantiated. This is equivalent to blocking until (1) the first *and* second arguments are both instantiated *or* (2) the third argument is instantiated. A reversal of `merge` would genuinely be useful because permutation generation frequently arises in solving constraint satisfaction problems. For example, figure 3 illustrates how a magic square program might be constructed in terms of `msort`. The cells of the square are constrained to be unique digits. The delay declaration on `sum` asserts that it cannot be invoked until its first, second and third arguments are all instantiated to non-variable terms.

The attempt at reversing mergesort given in figure 2 is not completely naïve. For instance, `split` will terminate if any of its arguments are bound to a rigid list (a list of predetermined length whose elements are not necessarily fixed). This is reflected in the delay declaration which asserts that `split` should block until any of its arguments are bound to a non-variable term. In fact a query such as `msort(L, [1, 2, 3])` *does* enumerate all its solutions; the problem is that it loops after doing so. This problem stems from the fact that `msort([X|L1], S1)` and `msort([Y|L2], S2)` can further instantiate `L1` and `L2`, which in turn enables `split(Xs, L1, L2)` to be activated to instantiate `Xs` to a list of 2 elements, so that `L` is of length 4. This process can continue without bound.

The mergesort program can be transformed into a version where termination is guaranteed for all goals. In particular, for a goal of the form `msort(L1, L2)`, where `L1` or `L2` are lists of numbers, then the computation cannot reduce to a state which contains a sub-goal that suspends indefinitely. More-



```

magic(Square) :-
  Square = [A, B, C, D, E, F, G, H, I],
  sum(A, B, C, Sum), sum(D, E, F, Sum), sum(G, H, I, Sum), sum(A, D, G, Sum),
  sum(B, E, H, Sum), sum(C, F, I, Sum), sum(A, E, I, Sum), sum(G, E, C, Sum),
  msort([- | Square], [0, 1, 2, 3, 4, 5, 6, 7, 8, 9]).
:- block sum(-, ?, ?, ?), sum(?, -, ?, ?), sum(?, ?, -, ?).
sum(A, B, C, Sum) :- Sum is A + B + C.

```

Figure 3. Magic square generation using mergesort

over, if it succeeds, then the set of answers produced is complete with respect to the declarative semantics. The transformed program is shown in figure 4. Each recursive predicate in the original program, namely `msort`, `split` and `merge`, is augmented by another predicate, namely `msort_depth`, `split_depth` and `merge_depth`, which computes a depth bound for that predicate. In addition, each predicate is replaced with an instrumented version, distinguished with the `sdr` suffix, which keeps track of the depth of the recursive calls. The auxiliary predicates that support the transformation are given separately in figure 5.

More precisely, the predicate `msort_depth(L1, L2, D)` calculates the lengths of the lists `L1` and `L2` with the predicate `list_length`. The `list_length(L, S, D)` predicate traverses the spine of the list `L`, blocking whenever a variable is encountered. Once `L` is rigid, `S` is bound to the length of `L`. This predicate uses an accumulator so that no more one pass is required over `L`. In fact, if `D` is instantiated before the pass is completed, then the pass can be aborted (see the penultimate paragraph of this section).

Once `L1` or `L2` are found to be rigid, the auxiliary predicate `msort_aux` calculates an upper bound `D` on the number of times that `msort_sdr` can call itself. Because of its divide-and-conquer nature, `msort` can call itself precisely  $\lceil \log_2(l) \rceil$  times where  $l$  is the length of `L1`, or equivalently, `L2`. Once `D` is instantiated to this value, the call to `msort_sdr(L, S, D)` can proceed. The purpose of the last argument is to ensure finiteness of the subsequent computation. The key point is that `D` is an upper bound on the number of recursive calls to `msort` in any successful derivation. Thus, by failing any derivation that exceeds this bound, termination is guaranteed without losing completeness. Observe that the depth bound for `msort` is only required to bound the depth of the predicates within the SCC for `msort`. Indeed,  $\lceil \log_2(l) \rceil$  does not bound the depth of the sub-computations for `split` and `merge`. Thus depth bounds can be calculated in an entirely modular fashion by considering each SCC in isolation. Note too that the delay declarations supplied in figure 2 are not relevant to either the transformation or the calculation of the depth bound.

For `split`, a depth bound on the number of recursive calls can be calculated as soon as either the first, second or third argument is rigid. This is because the depth bound is only required to be an *upper* bound. For instance, observe that if `split` is called with the third argument of length 2, and the call succeeds, then the first argument must be of length 4 or 5. Hence the number of recursive calls is either 4 or 5. Either way,  $2l + 1$  is an upper bound on the depth where  $l$  is the length of the third argument. Similarly,  $2l$  is an upper bound on the depth where  $l$  is the length of the second argument. The calculation of an exact depth bound would require the rigidity of both the second and third arguments, which would unnecessarily impede the execution of `split`. In the case of `merge`, a depth bound can be calculated as soon as either the first two arguments are rigid, or the third argument is rigid, but not earlier. In the case of `msort`, less precise upper bound could be applied in `msort_aux` to save distinguishing two cases by modifying the second clause to `D is integer(ceiling(log(D1 + 1)/log(2)))`, which is sufficient for

```

msort(L, S) :- msort_depth(L, S, D), msort_sdr(L, S, D).
:- block msort_sdr(?, ?, -).
msort_sdr([], [], -).
msort_sdr([X], [X], -).
msort_sdr([X, Y | Xs], S, D) :- 1 ≤ D, D1 is D - 1,
    split(Xs, L1, L2), msort_sdr([X|L1], S1, D1), msort_sdr([X|L2], S2, D1), merge(S1, S2, S).
split(L, L1, L2) :- split_depth(L, L1, L2, D), split_sdr(L, L1, L2, D).
:- block split_sdr(?, ?, ?, -).
split_sdr([], [], [], -).
split_sdr([X | L], [X | L1], L2, D) :- 1 ≤ D, D1 is D - 1, split_sdr(L, L2, L1, D1).
merge(L1, L2, L) :- merge_depth(L1, L2, L, D), merge_sdr(L1, L2, L, D).
:- block merge_sdr(?, ?, ?, -).
merge_sdr([], Ys, Ys, -).
merge_sdr(Xs, [], Xs, -).
merge_sdr([X|Xs], [Y|Ys], [X|Zs], D) :- 1 ≤ D, D1 is D - 1,
    leq(X, Y), merge_sdr(Xs, [Y|Ys], Zs, D1).
merge_sdr([X|Xs], [Y|Ys], [Y|Zs], D) :- 1 ≤ D, D1 is D - 1,
    gt(X, Y), merge_sdr([X|Xs], Ys, Zs, D1).

```

Figure 4. Control generation transformation applied to the mergesort program

avoiding a problematic zero logarithm. Alternatively putting  $D = D1$  would avoid a logarithm altogether. This illustrates a trade-off latent in the transformation: the tighter the bound, the earlier an otherwise non-terminating derivation is failed; yet this advantage has to be weighted against the disadvantages of manipulating a tighter bound. The extent to which bounds can be derived automatically (and thus the extent to which the transformation can be automated) depends on the particular form of the bound. Bounds that can be expressed as systems of linear inequations, for example, can be discovered by applying program analysis techniques in conjunction with program transformation techniques (as is explained in section 6.1). The key point is that the transformation is parameterised by the depth bounds and possibility of substituting one bound with a more tractable, albeit weaker, bound provides a route to automation.

To ensure termination, the computation rule is required to select each test  $1 \leq D$  before any other body atom. Any system implementing left-to-right selection rule (with delay) will satisfy this requirement. Thus termination is enforced with delays without resorting to a local computation rule [27]. Observe too that `list_length` predicate has 3 rather than 2 arguments. This is to save traversing lists unnecessarily. In the predicate `msort_depth`, for example, if the call `list_length(L, D1, D)` instantiates  $D1$ , then `msort_aux(D1, D2, D)` will bind  $D$  so that the length of  $S$  is inconsequential. Hence `list_length` terminates immediately if its third argument is instantiated.

Interestingly, the transformation yields a completely reversible mergesort program that possesses unexpected, though useful, computational properties. For instance, the query `msort(L, [1, 2, 2])` produces the 3 solutions:  $L = [1, 2, 2]$ ,  $L = [2, 1, 2]$  and  $L = [2, 2, 1]$  without producing any solution multiply. The reversal of `msort` is thus a combination generator rather than a permutation generator. A combination

```

msort_depth(L, S, D) :- list_length(L, D1, D), list_length(S, D2, D), msort_aux(D1, D2, D).
split_depth(L, L1, L2, D) :-
    list_length(L, S, D), list_length(L1, S1, D), list_length(L2, S2, D), split_aux(S, S1, S2, D).
merge_depth(L1, L2, L, D) :-
    list_length(L1, S1, D), list_length(L2, S2, D), list_length(L, S, D), merge_aux(S1, S2, S, D).
:- block msort_aux(-, -, ?).
msort_aux(D1, D2, D) :- D1 = D2, D1 = 0, !, D = 0.
msort_aux(D1, D2, D) :- D1 = D2, D is integer(ceiling(log(D1) / log(2))).
:- block split_aux(-, -, -, ?).
split_aux(S, _, _, D) :- nonvar(S), !, D = S.
split_aux(_, S1, _, D) :- nonvar(S1), !, D is 2 * S1.
split_aux(_, _, S2, D) :- D is 2 * S2 + 1.
:- block merge_aux(-, ?, -, ?), merge_aux(?, -, -, ?).
merge_aux(S1, S2, _, D) :- nonvar(S1), nonvar(S2), !, D is S1 + S2 - 1.
merge_aux(_, _, S, D) :- D is S - 1.
list_length(L, S, D) :- list_length(L, 0, S, D).
:- block list_length(-, ?, ?, -).
list_length(_, _, _, D) :- nonvar(D), !.
list_length([], S, S, _).
list_length([_:L], Acc, S, D) :- Acc1 is Acc + 1, list_length(L, Acc1, S, D).

```

Figure 5. Auxiliary predicates supporting the transformation

generator would be difficult to program directly without some notion of state to filter out replicated solutions. Incidentally, with the version `msort` listed in figure 4, `magic(Square)` finds two solutions, namely `Square = [6, 1, 8, 7, 5, 3, 2, 9, 4]` and `Square = [3, 8, 1, 2, 4, 6, 7, 0, 5]` (the other solutions are reflections). In the first case the rows, columns and diagonals all sum to 15 and in the other case they all sum to 12.

## 4. Theoretical foundations

To provide a sound theoretical basis for termination of delay logic programs, it is natural to build on the preceding theoretical foundations established for logic programs. Only definite programs (programs that exclude negation) are considered in this study.

### 4.1. Level mappings, norms and boundedness

The fundamental idea underlying all termination proofs is to define an order on the goals that can occur within a derivation. Given a program  $P$  and goal  $G_0$ , the finiteness of derivation  $G_0, G_1, G_2, \dots$  is in principle straightforward to demonstrate: it is sufficient to construct a well-founded order  $<$  such that  $G_{i+1} < G_i$  for all  $i \geq 0$ . The problem is to find such an order. To simplify the problem, it is convenient to define the order on abstractions of goals rather than on the goals themselves. Thus the order  $<$  is

defined such that  $G' < G$  holds iff  $\mathcal{A}(G') < \mathcal{A}(G)$  holds where  $\mathcal{A}$  is an abstraction function. For example,  $\mathcal{A}$  might be defined to map each goal  $G$  to a multiset of natural numbers, where each atom in  $G$  maps to a single number in the multiset. The approach leads to the concept of a level mapping [11].

**Definition 4.1.** Let  $P$  be a program. A level mapping for  $P$  is a function  $|\cdot| : B_P \rightarrow \mathbb{N}$  from the Herbrand base of  $P$  to the set of natural numbers  $\mathbb{N}$ . For an atom  $A \in B_P$ ,  $|A|$  denotes the level of  $A$ .

The reader is referred to Lloyd [23] for the standard definitions of the Herbrand universe, Herbrand base, Herbrand interpretation, Herbrand model, clauses, substitutions, unifiers, most general unifiers, resolvents, derivations, *etc.* Since a level mapping is defined over the Herbrand base, it is not defined for non-ground atoms. However, the concept can be straightforwardly lifted [5].

**Definition 4.2.** An atom  $A$  is bounded wrt a level mapping  $|\cdot|$  if  $|\cdot|$  is bounded on the set  $[A]$  of variable free instances of  $A$ . If  $A$  is bounded then  $[[A]]$  denotes the maximum that  $|\cdot|$  takes on  $[A]$ .

The importance of boundedness cannot be over stressed. Since goals which are ground cannot be used to compute values, they are the exception rather than the norm in logic programming. Thus practical termination proofs must deal with non-ground goals and boundedness provides the basis for this.

**Example 4.1.** Let  $P$  be the program:

$p(a, X) :- p(b, X).$   
 $p(b, a).$   
 $p(b, b).$

The function  $|\cdot| : \{p(a, a), p(a, b), p(b, a), p(b, b)\} \rightarrow \mathbb{N}$  defined by  $|p(a, a)| = 34$ ,  $|p(a, b)| = 12$ ,  $|p(b, a)| = 0$  and  $|p(b, b)| = 27$  is a level mapping for  $P$ . The atom  $p(a, X)$  is bounded wrt  $|\cdot|$  since  $[p(a, X)] = \{p(a, a), p(a, b)\}$  and  $\max(\{|p(a, a)|, |p(a, b)|\}) = \max(\{34, 12\}) = 34$ .

Level mappings are usually defined in terms of norms. A norm is a mapping from terms to natural numbers which provides some measure of the size of a term.

**Example 4.2.** The list-length norm  $|\cdot|_{list-length} : U_P \rightarrow \mathbb{N}$  and the term-size norm  $|\cdot|_{term-size} : U_P \rightarrow \mathbb{N}$  from the Herbrand universe to the natural numbers can be defined by

$$|t|_{list-length} = \begin{cases} 1 + |t_2|_{list-length} & \text{if } t = [t_1|t_2] \\ 0 & \text{otherwise} \end{cases} \quad |f(t_1, \dots, t_n)|_{term-size} = 1 + \sum_{i=1}^n |t_i|_{term-size}$$

Then, for example,  $[[a, b, f(a)]]_{list-length} = 3$  and  $|f(a, g(b))|_{term-size} = 1 + |a|_{term-size} + |g(b)|_{term-size} = 1 + 1 + 1 + |b|_{term-size} = 4$ .

## 4.2. Atom selection

In the classic level mapping based approaches to termination [1, 5], a fundamental requirement is that only bounded atoms are selected. The reason is that, in general, when unbounded atoms are selected for resolution, it becomes more difficult to reason about the termination of the subsequent computation. This

approach can still be applied when considering flexible computation rules. Moreover, delay declarations provide a mechanism to control this directly by delaying atoms until they become bounded. This idea is formally captured in the following definition which is due to Marchiori and Teusink [26, 27].

**Definition 4.3.** A delay declaration for a predicate  $p$  is safe wrt a level mapping  $|\cdot|$  if for every atom  $A$  with predicate symbol  $p$ , if  $A$  satisfies its delay declaration, then  $A$  is bounded wrt  $|\cdot|$ .

### 4.3. Covers

To determine whether or not an atom is bounded when it is selected, requires a consideration of the atoms that have been (partially) resolved before the selection of the atom. The following definitions of cover, originally proposed by Marchiori and Teusink [26, 27], capture this idea. Note that in the following definition,  $vars(o)$  denotes the set of variables within the syntactic object  $o$ ;  $body(c)$  is the set of body atoms occurring in a clause  $c$ ; and  $dom(\theta)$  is the set of variables which constitute the domain of a substitution  $\theta$ .

**Definition 4.4.** Let  $c = H:-B_1, \dots, B_n$  be a clause and  $|\cdot|$  a level mapping. Let  $A \in body(c)$  and  $D \subseteq body(c)$  such that  $A \notin D$ . Then  $D$  is a direct cover for  $A$  wrt  $|\cdot|$  in  $c$ , if there exists a substitution  $\theta$  such that:

- $\theta(A)$  is bounded wrt  $|\cdot|$  and
- $dom(\theta) \subseteq vars(H) \cup vars(D)$ .

A direct cover  $D$  for  $A$  is minimal if no proper subset of  $D$  is a direct cover for  $A$ . The set of minimal direct covers of  $A$  wrt  $|\cdot|$  in  $c$  is denoted by  $mdcovers_{|\cdot|,c}(A)$ .

Intuitively, a direct cover of an atom  $A$  in a clause  $c$  is a subset  $D$  of the body atoms of  $c$  such that for some instantiation  $\theta$  of the variables in  $D$ ,  $\theta(A)$  is bounded. Note that a body atom may have zero, one or more (minimal) direct covers. In particular, an atom  $A$  will have no direct cover when, in order for  $A$  to become bounded, it is necessary to instantiate a variable of  $A$  which does not occur elsewhere in the clause. On the other hand, the atom  $A$  will have the empty set as its only minimal cover if  $A$  is bounded whenever the head of the clause is ground.

**Example 4.3.** Consider the quicksort program listed in figure 1 and the level mapping  $|\cdot|$  defined by:

$$|qsort(x, y)| = y' + 1 \quad |part(w, x, y, z)| = y' + z' \quad |append(x, y, z)| = z'$$

where  $y' = |y|_{list-length}$  and  $z' = |z|_{list-length}$ . Then

$$\begin{aligned} mdcovers_{|\cdot|,c}(part(Xs, X, L, B)) &= \{\{qsort(L, Ls), qsort(B, Bs)\}\} \\ mdcovers_{|\cdot|,c}(qsort(L, Ls)) &= \{\{append(Ls, [X|Bs], Ys)\}\} \\ mdcovers_{|\cdot|,c}(qsort(B, Bs)) &= \{\{append(Ls, [X|Bs], Ys)\}\} \\ mdcovers_{|\cdot|,c}(append(Ls, [X|Bs], Ys)) &= \{\emptyset\} \end{aligned}$$

where  $c$  is the second clause of `qsort`. Note that  $mdcovers_{|\cdot|,c}(append(Ls, [X|Bs], Ys)) = \{\emptyset\}$  because  $Ys \in vars(qsort([X|Xs], Ys))$ . In this example each atom has exactly one minimal direct cover.

Marchiori and Teusink [27] explain the idea of cover in terms of “potential activators”. Suppose that each predicate of the quicksort program is equipped with a delay declaration that is safe wrt the level mapping  $|\cdot|$  of example 4.3. Then the goal  $\text{qsort}(L, Ls)$  can only be selected when  $\text{qsort}(L, Ls)$  is bounded wrt  $|\cdot|$ , that is, when  $Ls$  is instantiated to a rigid list. Resolution of  $\text{append}(Ls, [X|Bs], Ys)$  can potentially activate this goal by instantiating  $Ls$  thereby making  $\text{qsort}(L, Ls)$  selectable. The (minimal) direct cover codifies this potential producer-consumer relationship among the body atoms of a clause. It can also express multiple producer-single consumer relationships. The goal  $\text{part}(Xs, X, L, B)$  is only selectable if it is bounded wrt  $|\cdot|$ , or equivalently, when  $L$  and  $B$  are instantiated to lists of fixed length. The goals  $\text{qsort}(L, Ls)$  and  $\text{qsort}(B, Bs)$  are potential activators for  $\text{part}(Xs, X, L, B)$  since, when combined, they can potentially instantiate  $L$  and  $B$  to rigid lists. Note, however, that although (minimal) direct cover encapsulates the idea of potential activation, it misses the idea that activators themselves need to be activated. This iterative application of potential activation leads to the concept of cover.

**Definition 4.5.** Let  $c = H:-B_1, \dots, B_n$  be a clause and  $|\cdot|$  a level mapping. Let  $A \in \text{body}(c)$  and  $C \subseteq \text{body}(c)$  such that  $A \notin C$ . Then  $C$  is a *cover* for  $A$  wrt  $|\cdot|$  in  $c$ , if  $\langle A, C \rangle$  is an element of the least set  $S$  such that:

- $\langle A, \emptyset \rangle \in S$  if  $\emptyset \in \text{mdcovers}_{|\cdot|,c}(A)$  or
- $\langle A, C \rangle \in S$  if  $A \notin C$ ,  $\{A_1, \dots, A_k\} \in \text{mdcovers}_{|\cdot|,c}(A)$ ,  $\langle A_i, C_i \rangle \in S$  for all  $i \in [1, k]$  and  $C = \cup_{i=1}^k \{A_i\} \cup C_i$ .

The set of covers of  $A$  wrt  $|\cdot|$  in  $c$  is denoted by  $\text{covers}_{|\cdot|,c}(A)$ .

Interestingly, although the cover relation is a kind of closure of the direct cover relation, it is not a transitive one; a direct cover of an atom is not necessarily a cover of that atom. This is illustrated by  $\text{part}(Xs, X, L, B)$ . The minimum direct cover specifies that this goal can only be activated by  $\text{qsort}(L, Ls)$  and  $\text{qsort}(B, Bs)$ . However, these can only be activated by  $\text{append}(Ls, [X|Bs], Ys)$ . Hence the two  $\text{qsort}(L, Ls)$  and  $\text{qsort}(B, Bs)$  do not qualify as cover for  $\text{part}(Xs, X, L, B)$ ;  $\text{part}(Xs, X, L, B)$  can only be activated once  $\text{append}(Ls, [X|Bs], Ys)$  has activated both  $\text{qsort}(L, Ls)$  and  $\text{qsort}(B, Bs)$ .

**Example 4.4.** Consider the program `qsort`, the level mapping  $|\cdot|$  and clause  $c$  of example 4.3, and the set  $S$  as defined in definition 4.5. Then

- Since  $\emptyset \in \text{mdcovers}_{|\cdot|,c}(\text{append}(Ls, [X|Bs], Ys))$  it follows  $\langle \text{append}(Ls, [X|Bs], Ys), \emptyset \rangle \in S$ ;
- Since  $\{\text{append}(Ls, [X|Bs], Ys)\} \in \text{mdcovers}_{|\cdot|,c}(\text{qsort}(L, Ls))$  and  $\langle \text{append}(Ls, [X|Bs], Ys), \emptyset \rangle \in S$  then  $\langle \text{qsort}(L, Ls), C \rangle \in S$  where  $C = \{\text{append}(Ls, [X|Bs], Ys)\} \cup \emptyset$ ;
- It follows similarly that  $\langle \text{qsort}(B, Bs), \{\text{append}(Ls, [X|Bs], Ys)\} \rangle \in S$ ;
- Finally  $\langle \text{part}(Xs, X, L, B), C \rangle \in S$  where  $C = (\{\text{qsort}(L, Ls)\} \cup \{\text{append}(Ls, [X|Bs], Ys)\}) \cup (\{\text{qsort}(B, Bs)\} \cup \{\text{append}(Ls, [X|Bs], Ys)\})$  because
  - $\{\text{qsort}(L, Ls), \text{qsort}(B, Bs)\} \in \text{mdcovers}_{|\cdot|,c}(\text{part}(Xs, X, L, B))$  and
  - $\langle \text{qsort}(L, Ls), \{\text{append}(Ls, [X|Bs], Ys)\} \rangle \in S$  and
  - $\langle \text{qsort}(B, Bs), \{\text{append}(Ls, [X|Bs], Ys)\} \rangle \in S$ .

$$\begin{aligned}
 \text{Hence } \text{covers}_{|\cdot|,c}(\text{part}(Xs, X, L, B)) &= \{\{\text{qsort}(L, Ls), \text{qsort}(B, Bs), \text{append}(Ls, [X|Bs], Ys)\}\} \\
 \text{covers}_{|\cdot|,c}(\text{qsort}(L, Ls)) &= \{\{\text{append}(Ls, [X|Bs], Ys)\}\} \\
 \text{covers}_{|\cdot|,c}(\text{qsort}(B, Bs)) &= \{\{\text{append}(Ls, [X|Bs], Ys)\}\} \\
 \text{covers}_{|\cdot|,c}(\text{append}(Ls, [X|Bs], Ys)) &= \{\emptyset\}
 \end{aligned}$$

#### 4.4. Delay recurrency

Using the notion of cover, Marchiori and Teusink [26, 27] introduced the class of delay recurrent programs. It was intended that programs lying within this class would be terminating under a dynamic selection rule.

**Definition 4.6.** Let  $P$  be a program,  $|\cdot|$  a level mapping and  $I$  an interpretation for  $P$ . A clause  $c = H:-B_1, \dots, B_n$  is delay recurrent wrt  $|\cdot|$  and  $I$  iff

- $I$  is a model for  $c$  and
- for all  $i \in [1, n]$ , for every cover  $C \in \text{covers}_{|\cdot|,c}(B_i)$  and for every grounding substitution  $\theta$  for  $c$  such that  $I \models \theta(C)$ , it follows that  $|\theta(H)| > |\theta(B_i)|$ .

A program  $P$  is delay recurrent wrt  $|\cdot|$  and  $I$  iff every clause of  $P$  is delay recurrent wrt  $|\cdot|$  and  $I$ .

To be precise, definition 4.6 differs from that in [26, 27] since the original contains some slight redundancy/ambiguity. The above definition, however, reflects the intentions of Marchiori and Teusink [25].

**Example 4.5.** Let  $|\cdot|$  be the level mapping of example 4.3 and  $I$  the interpretation

$$\begin{aligned}
 \{\text{qsort}(x, y) \mid |x|_{\text{list-length}} &= |y|_{\text{list-length}}\} \cup \\
 \{\text{part}(x, w, y, z) \mid |x|_{\text{list-length}} &= |y|_{\text{list-length}} + |z|_{\text{list-length}}\} \cup \\
 \{\text{append}(x, y, z) \mid |z|_{\text{list-length}} &= |x|_{\text{list-length}} + |y|_{\text{list-length}}\}
 \end{aligned}$$

and consider the second clause of **qsort** dubbed  $c$ . To check that  $I$  is a model for  $c$  let  $\theta = \{X \mapsto t_1, Xs \mapsto t_2, Ys \mapsto t_3, L \mapsto t_4, B \mapsto t_5, Ls \mapsto t_6, Bs \mapsto t_7\}$  be a grounding substitution for  $c$  such that  $I \models \{\text{part}(t_2, t_1, t_4, t_5), \text{qsort}(t_4, t_6), \text{qsort}(t_5, t_7), \text{append}(t_6, [t_1|t_7], t_3)\}$ .  $I$  is a model for  $c$  because

$$\begin{aligned}
 |[t_1|t_2]|_{\text{list-length}} &= 1 + |t_2|_{\text{list-length}} = 1 + (|t_4|_{\text{list-length}} + |t_5|_{\text{list-length}}) \\
 &= 1 + (|t_6|_{\text{list-length}} + |t_7|_{\text{list-length}}) \\
 &= |t_6|_{\text{list-length}} + |[t_1|t_7]|_{\text{list-length}} = |t_3|_{\text{list-length}}
 \end{aligned}$$

Now consider the body atom  $\text{part}(Xs, X, L, B)$  and let  $C \in \text{covers}_{|\cdot|,c}(\text{part}(Xs, X, L, B))$ . It follows that  $C = \{\text{qsort}(L, Ls), \text{qsort}(B, Bs), \text{append}(Ls, [X|Bs], Ys)\}$ . Consider again  $\theta$  as it is a grounding substitution for  $c$  and suppose  $I \models \{\text{qsort}(t_4, t_6), \text{qsort}(t_5, t_7), \text{append}(t_6, [t_1|t_7], t_3)\}$ . Then

$$\begin{aligned}
 |\text{qsort}([t_1|t_2], t_3)| &= |t_3|_{\text{list-length}} + 1 = (|t_6|_{\text{list-length}} + |t_7|_{\text{list-length}} + 1) + 1 \\
 &= (|t_4|_{\text{list-length}} + |t_5|_{\text{list-length}} + 1) + 1 \\
 &> |t_4|_{\text{list-length}} + |t_5|_{\text{list-length}} = |\text{part}(t_2, t_1, t_4, t_5)|
 \end{aligned}$$

It is straightforward to check that condition 2 of definition 4.6 holds for every other body atom of  $c$ . Hence  $c$  is delay recurrent wrt  $|\cdot|$  and  $I$ .

The intention behind the definition of delay recurrency is that a delay recurrent program  $P$ , when augmented with a set of safe delay declarations for the predicates of  $P$ , only admits finite derivations. The delay declarations handle the local boundedness issue, but there is still the global boundedness issue to consider. Suppose  $C$  is a cover for an atom  $B$  in a delay recurrent clause, and  $\theta$  is a partial answer substitution for  $C$  such that  $\theta(B)$  is bounded (note that  $\theta$  may not necessarily be a *correct* answer substitution since the atoms in  $C$  have not yet been fully resolved). At this point  $\theta$  speculatively binds the variables of  $B$  since it is not yet known whether or not there exists some substitution  $\sigma$  such that  $I \models \sigma(\theta(C))$ . If  $\theta(B)$  is selected at this point an infinite computation may arise since there is no guarantee that the level of the head is greater than the level of  $\theta(B)$ . Instead, by fully resolving each atom in  $C$  such that a correct answer substitution  $\theta$  is obtained,  $\theta(B)$  can be selected safely since  $I \models \sigma(\theta(C))$  for all  $\sigma$ , whence by condition 2 of delay recurrency, the level of  $\theta(B)$  is less than the level of the head. Full resolution of  $C$  can be achieved by using a local selection rule. To reiterate, a local selection rule only selects the most recently introduced atoms in a derivation and thus completely resolves sub-computations before proceeding with the main computation [43]. The notion is formally defined below.

**Definition 4.7.** For a goal  $G = A_1, \dots, A_m$  the  $i$ th atom of  $G$  is  $A_i$ . Let  $G_0, G_1, G_2, \dots, G_k, \dots$  be a derivation. The age of the  $i$ th atom in  $G_k$ , denoted  $age_{G_k}(i)$ , is defined as follows:

- If  $G_0 = A_1, \dots, A_m$ , then  $age_{G_0}(i) = 0$  for all  $i \in [1, m]$ ;
- If  $G_k = A_1, \dots, A_m$  and  $G_{k+1} = \theta(A_1, \dots, A_{s-1}, B_1, \dots, B_n, A_{s+1}, \dots, A_m)$ , then

$$age_{G_{k+1}}(i) = \begin{cases} age_{G_k}(i) + 1 & \text{for all } i \in [1, s-1] \\ 0 & \text{for all } i \in [s, s+n-1] \\ age_{G_k}(i-n+1) + 1 & \text{for all } i \in [s+n, n+m-1] \end{cases}$$

For a goal  $G = A_1, \dots, A_m$ , an atom  $A_i$  is introduced in  $G$  if  $age_G(A) = 0$ . An atom  $A_i$  is most recently introduced in  $G$  iff  $age_G(i) \leq age_G(j)$  for all  $j \in [1, m]$ .

**Definition 4.8.** A derivation  $G_0, G_1, G_2, \dots, G_k, \dots$  applies local selection iff whenever an atom  $A_i$  is selected for resolution from  $G_k$  then  $A_i$  is most recently introduced in  $G_k$ .

**Definition 4.9.** A derivation  $G_0, G_1, G_2, \dots, G_k, \dots$  satisfies the delay declarations of  $P$  iff whenever an atom  $A_i$  is selected for resolution from  $G_k$  then  $A_i$  satisfies the delay declarations of  $P$ .

The main result regarding delay recurrent programs [26, 27] can now be stated:

**Theorem 4.1.** Let  $P$  be a program with a delay declaration for each predicate in  $P$ . Let  $|\cdot|$  be a level mapping and  $I$  an interpretation. Suppose that

- $P$  is delay recurrent wrt  $|\cdot|$  and  $I$ , and
- the delay declarations for  $P$  are safe wrt  $|\cdot|$

Then every derivation that satisfies the delay declarations of  $P$ , whilst applying local selection, is finite.



#### 4.5. Semi-delay recurrency

Marchiori and Teusink [26, 27] noticed that boundedness of atoms could be enforced by using delay declarations but did not fully exploit this fact. Their definition requires a decrease in the level mapping from the head to the non-recursive body atoms, when in fact the boundedness of selected atoms is already guaranteed by safe delay declarations. To relax this restriction, notion of predicate dependency is introduced to distinguish between recursive and non-recursive body atoms.

**Definition 4.10.** Let  $p, q \in \Pi$  where  $\Pi$  is the set of predicate symbols in a logic program  $P$ . Then  $p$  directly-depends-on  $q$  iff  $p(t_1, \dots, t_k) :- B_1, \dots, B_n \in P$  and  $B_i = q(s_1, \dots, s_l)$  for some  $i \in [1, n]$ . The depends-on relation, denoted  $\sqsupseteq$ , is defined as the reflexive, transitive closure of the directly-depends-on relation. If  $p \sqsupseteq q$  and  $q \sqsupseteq p$  then  $p$  and  $q$  are mutually dependent and this is denoted by  $p \simeq q$ .

**Example 4.6.** Returning to the quicksort program listed in figure 1,  $\text{qsort} \sqsupseteq \text{part}$ ,  $\text{qsort} \sqsupseteq \text{append}$ ,  $\text{part} \sqsupseteq \text{leq}$ ,  $\text{part} \sqsupseteq \text{gt}$ ,  $\text{qsort} \sqsupseteq \text{leq}$  and  $\text{qsort} \sqsupseteq \text{gt}$ . For this program,  $p \simeq q$  iff  $p = q$ .

The following new definition will prove useful for defining a large class of terminating programs which permit coroutining. In what follows  $\text{rel}(A)$  denotes the predicate symbol of the atom  $A$ .

**Definition 4.11.** Let  $|\cdot|$  be a level mapping and  $I$  an interpretation for a program  $P$ . A clause  $c = H :- B_1, \dots, B_n$  is semi-delay recurrent wrt  $|\cdot|$  and  $I$  iff

- $I$  is a model for  $c$  and
- for all  $i \in [1, n]$  such that  $\text{rel}(H) \simeq \text{rel}(B_i)$ , for every cover  $C \in \text{covers}_{|\cdot|,c}(B_i)$  and for every grounding substitution  $\theta$  for  $c$  such that  $I \models \theta(C)$ , it follows that  $|\theta(H)| > |\theta(B_i)|$ .

A program  $P$  is semi-delay recurrent wrt  $|\cdot|$  and  $I$  iff every clause of  $P$  is semi-delay recurrent wrt  $|\cdot|$  and  $I$ .

Observe in this definition that there is no restriction placed on the relation between the level of the head of the clause and the level of the non-recursive body atoms.

**Example 4.7.** Let  $I$  be the interpretation and  $|\cdot|$  the level mapping of example 4.5. As before,  $I$  is a model for the second clause  $c$  of the  $\text{qsort}$  predicate. Now  $\text{rel}(\text{qsort}([X|Xs], Ys)) \simeq \text{rel}(\text{qsort}(L, Ls))$  and  $\text{rel}(\text{qsort}([X|Xs], Ys)) \simeq \text{rel}(\text{qsort}(B, Bs))$ . Then

$$\begin{aligned} \text{covers}_{|\cdot|,c}(\text{qsort}(L, Ls)) &= \{\{\text{append}(Ls, [X|Bs], Ys)\}\} \\ \text{covers}_{|\cdot|,c}(\text{qsort}(B, Bs)) &= \{\{\text{append}(Ls, [X|Bs], Ys)\}\} \end{aligned}$$

Let  $\theta = \{X \mapsto t_1, Xs \mapsto t_2, Ys \mapsto t_3, L \mapsto t_4, B \mapsto t_5, Ls \mapsto t_6, Bs \mapsto t_7\}$  be a grounding substitution for  $c$  such that  $I \models \{\text{append}(t_6, [t_1|t_7], t_3)\}$ . Then

$$\begin{aligned} |\text{qsort}([t_1|t_2], t_3)| &= |t_3|_{\text{list-length}} + 1 \\ &= (|t_6|_{\text{list-length}} + |t_7|_{\text{list-length}} + 1) + 1 \\ &> |t_6|_{\text{list-length}} + 1 \\ &= |\text{qsort}(t_4, t_6)| \end{aligned} \qquad \begin{aligned} |\text{qsort}([t_1|t_2], t_3)| &= |t_3|_{\text{list-length}} + 1 \\ &> |t_7|_{\text{list-length}} + 1 \\ &= |\text{qsort}(t_5, t_7)| \end{aligned}$$

Hence  $c$  is semi-delay recurrent wrt  $|\cdot|$  and  $I$ .

It would be straightforward to prove that theorem 4.1 still holds if the program is replaced by one which is semi-delay recurrent, but a more significant result can be obtained. Observe that a local selection rule is used to ensure that a cover of an atom is completely resolved before the atom itself is selected. Notice, however, that for semi-delay recurrency, it is only necessary for the covers of the mutually recursive atoms to be resolved completely. This means that, following the resolution of these covers, an arbitrary amount of coroutining may take place amongst the remaining atoms of the clause.

To formalise a selection rule based on this idea, the notion of covering is lifted from the clause level to the goal level. A covering of a recursive atom  $A$  in a goal  $G$ , is the set of atoms in  $G$  which have yet to be resolved before  $A$  can be safely selected. An atom  $A$  may have more than one covering, though it is only necessary to fully resolve the atoms of one of them before the selection of  $A$ . Coverings of atoms will change during the course of a derivation as new atoms are introduced and others are fully resolved.

**Definition 4.12.** Let  $G_0, G_1, G_2, \dots$  be a derivation and  $|\cdot|$  a level mapping. A covering  $C$  for an atom  $A$  in a goal  $G_l$  wrt  $|\cdot|$ , denoted  $C \in \text{covers}_{|\cdot|, G_l}(A)$ , is defined as follows:

- Suppose  $G_0 = A_1, \dots, A_m$ . Then  $\emptyset \in \text{covers}_{|\cdot|, G_0}(A_i)$  for all  $i \in [1, m]$ .
- Suppose  $G_{l+1} = \theta(A_1, \dots, A_{s-1}, B_1, \dots, B_n, A_{s+1}, \dots, A_m)$  is the resolvent derived from  $G_l = A_1, \dots, A_s, \dots, A_m$  and  $c = H:-B_1, \dots, B_n$ , where  $A_s$  is the selected atom in  $G_l$  and  $\theta \in \text{mgu}(H, A_s)$ . Then
  - for all  $i \in [1, n]$ ,
    - \* if  $\text{rel}(H) \simeq \text{rel}(B_i)$  and  $C \in \text{covers}_{|\cdot|, c}(B_i)$  then  $\theta(C) \in \text{covers}_{|\cdot|, G_{l+1}}(\theta(B_i))$ ;
    - \* if  $\text{rel}(H) \not\simeq \text{rel}(B_i)$  then  $\emptyset \in \text{covers}_{|\cdot|, G_{l+1}}(\theta(B_i))$ .
  - for all  $i \in [1, m]$ , if  $i \neq s$ ,  $C \subseteq \{A_1, \dots, A_{i-1}, A_{i+1}, \dots, A_m\}$  and  $C \in \text{covers}_{|\cdot|, G_l}(A_i)$  then
    - \* if  $A_s \in C$  then  $\theta(C) \setminus \{\theta(A_s)\} \cup \theta(\{B_1, \dots, B_n\}) \in \text{covers}_{|\cdot|, G_{l+1}}(\theta(A_i))$ ;
    - \* if  $A_s \notin C$  then  $\theta(C) \in \text{covers}_{|\cdot|, G_{l+1}}(\theta(A_i))$ .

**Definition 4.13.** Let  $|\cdot|$  be a level mapping. A derivation  $G_0, G_1, G_2, \dots, G_k, \dots$  applies semi-local selection wrt  $|\cdot|$  iff whenever an atom  $A_i$  is selected for resolution from  $G_k$  then  $\emptyset \in \text{covers}_{|\cdot|, G_k}(A_i)$ .

Observe that, if an atom is selectable under a local selection rule, then it is selectable under a semi-local selection rule. The main result can now be stated:

**Theorem 4.2.** Let  $P$  be a program with a delay declaration for each predicate in  $P$ . Let  $|\cdot|$  be a level mapping and  $I$  an interpretation. Suppose that

- $P$  is semi-delay recurrent wrt  $|\cdot|$  and  $I$
- The delay declarations for  $P$  are safe wrt  $|\cdot|$

Then every derivation that satisfies the delay declarations of  $P$ , whilst applying semi-local selection wrt  $|\cdot|$ , is finite.

The termination argument is constructed in terms of a multiset ordering that is specific to the program. The multiset ordering  $\prec_{mul}$  is itself formulated in terms of an ordering  $\prec$  on  $\Pi \times \mathbb{N}$ . The binary relation  $\prec$  is defined  $\langle p, m \rangle \prec \langle q, n \rangle$  iff  $(p \simeq q \text{ and } m < n)$  or  $(p \not\simeq q \text{ and } q \sqsupseteq p)$ . In the following definition, the brackets  $\{\!\{$  and  $\}\!\}$  denote the delimiters of a multiset, for example,  $\{\!\{ \langle p, m \rangle \}\!\}$  denotes the multiset that only contains the single pair  $\langle p, m \rangle$ . The multiset ordering is defined by  $s_1 \prec_{mul} s_2$  iff there exists  $\langle p_1, m_1 \rangle, \dots, \langle p_k, m_k \rangle \in s_1$  and  $\langle p, m \rangle \in s_2$  such that  $s_1 = (s_2 / \{\!\{ \langle p, m \rangle \}\!\}) \cup \{\!\{ \langle p_1, m_1 \rangle, \dots, \langle p_k, m_k \rangle \}\!\}$  and  $\langle p_i, m_i \rangle \prec \langle p, m \rangle$  for all  $i \in [1, k]$ . Since  $\Pi$  is finite and  $\mathbb{N}$  is well-ordered, it follows that  $\prec$  is well-ordered, hence  $\prec_{mul}$  is well-ordered [41]. It then remains to specify a mapping from each goal of a derivation to multiset such that a decrease can be observed between multisets of consecutive goals. Such a mapping is given below:

**Definition 4.14.** Let  $G_0, G_1, G_2, \dots$  be a derivation,  $|\cdot|$  be a level mapping,  $I$  an interpretation and  $G_k = A_1, \dots, A_n$  be a goal. For all  $i \in [1, n]$  define

$$|[G_k]_I^i| = \left\{ \begin{array}{l} \theta \text{ is a grounding substitution for } G_k \quad \wedge \\ |\theta(A_i)| + 1 \quad \left| \begin{array}{l} C \in \text{covers}_{|\cdot|, G_k}(A_i) \\ I \models \theta(C) \end{array} \right. \quad \wedge \end{array} \right\}$$

Furthermore, if the set  $|[G]_I^i|$  is finite for each  $i \in [1, n]$  then define

$$|[G]_I| = \{\!\{ \langle \text{rel}(A_1), \max |[G]_I^1| \rangle, \dots, \langle \text{rel}(A_n), \max |[G]_I^n| \rangle \}\!\}$$

where  $\max \emptyset = 0$  and  $\max \{n_1, \dots, n_k\} = n_i$  where  $n_i \geq n_j$  for all  $j \in [1, k]$ . (Note that the multiset  $|[G]_I|$  is well-defined if and only if each set  $|[G]_I^i|$  is finite.)

**Lemma 4.1.** Let  $|\cdot|$  be a level mapping and let  $G_0, G_1, G_2, \dots, G_k, G_{k+1}$  be a (partial) derivation that satisfies the delay declarations of  $P$ , whilst applying semi-local selection wrt  $|\cdot|$ . Let  $I$  be an interpretation for program  $P$  and suppose  $P$  is semi-delay recurrent wrt  $|\cdot|$  and  $I$ . Then if  $|[G_k]_I|$  is well-defined it follows that  $|[G_{k+1}]_I|$  is well-defined and  $|[G_{k+1}]_I| \prec_{mul} |[G_k]_I|$ .

**Proof:**

For brevity, let  $G = G_k$  and  $G' = G_{k+1}$ . Then  $G' = \theta(A_1, \dots, A_{s-1}, B_1, \dots, B_n, A_{s+1}, \dots, A_m)$  is the resolvent derived from  $G = A_1, \dots, A_s, \dots, A_m$  and  $c = H: -B_1, \dots, B_n$  where  $A_s$  is the selected atom in  $G$  and  $\theta \in \text{mgu}(A_s, H)$ . Suppose that  $|[G_k]_I|$  is well-defined.

- To show  $\langle \text{rel}(\theta(B_i)), \max |[G']_I^{s+i-1}| \rangle \prec \langle \text{rel}(A_s), \max |[G]_I^s| \rangle$  for all  $i \in [1, n]$ . Let  $i \in [1, n]$ . If  $\text{rel}(B_i) \not\simeq \text{rel}(H) \simeq \text{rel}(A_s)$  then  $\text{rel}(B_i) \sqsupseteq \text{rel}(A_s)$  and the result follows. Now suppose  $\text{rel}(B_i) \simeq \text{rel}(H)$ . By definition 4.13,  $\emptyset \in \text{covers}_{|\cdot|, G}(A_s)$ . By definition 4.14 it follows that  $\max |[G]_I^s| = \max \{|\sigma(A_s)| + 1 \mid \sigma \text{ is a grounding substitution for } G\}$ . Since the derivation satisfies the delay declarations of  $P$ ,  $A_s$  is bounded wrt  $|\cdot|$ , hence  $\max |[G]_I^s| = |[A_s]| + 1$ , whence  $|[A_s]| \geq |\theta(A_s)| = |\theta(H)|$ . Thus  $\max |[G]_I^s| > |\theta(H)|$ . By definition 4.14 it follows that

$$\begin{aligned} |[G']_I^{s+i-1}| &= |\theta(A_1, \dots, A_{s-1}, B_1, \dots, B_n, A_{s+1}, \dots, A_m)|_I^{s+i-1} \\ &= \left\{ \begin{array}{l} \sigma \text{ is a grounding substitution for } G' \quad \wedge \\ |\sigma(\theta(B_i))| + 1 \quad \left| \begin{array}{l} D \in \text{covers}_{|\cdot|, G'}(\theta(B_i)) \\ I \models \sigma(D) \end{array} \right. \quad \wedge \end{array} \right\} \end{aligned}$$

- Suppose  $[[G']_I^{s+i-1}] \neq \emptyset$ . Then there exists (1) a grounding substitution  $\sigma$  for  $G'$  and (2) a cover  $D \in \text{covers}_{|,G'}(\theta(B_i))$  such that (3)  $I \models \sigma(D)$  and (4)  $|\sigma(\theta(B_i))| + 1 = \max [[G']_I^{s+i-1}]$ . From (2) it follows that there exists a cover  $C \in \text{covers}_{|,c}(B_i)$  such that  $D = \theta(C)$ , hence by (3),  $I \models \sigma(\theta(C))$ . By (1),  $\sigma \circ \theta$  is a grounding substitution for  $B_1, \dots, B_n$ . Let  $\kappa$  be a grounding substitution for  $\sigma(\theta(H))$ . Then  $\kappa \circ \sigma \circ \theta$  is a grounding substitution for  $c$  such that  $I \models \kappa(\sigma(\theta(C)))$  since  $\kappa(\sigma(\theta(C))) = \sigma(\theta(C))$ . Hence, by the semi-delay recurrency of  $c$ ,  $|\kappa(\sigma(\theta(H)))| > |\kappa(\sigma(\theta(B_i)))| = |\sigma(\theta(B_i))|$ . Thus  $\max [[G']_I^{s+i-1}] \leq |\kappa(\sigma(\theta(H)))| = |\kappa(\sigma(\theta(H)))| \leq |\theta(H)|$ . It follows  $\max [[G']_I^s] > \max |G'|_I^{s+i-1}$ , hence  $|G'|_I^{s+i-1}$  is finite.
- Suppose  $[[G']_I^{s+i-1}] = \emptyset$ . Then  $[[G']_I^{s+i-1}]$  is trivially finite. Moreover  $\max [[G']_I^{s+i-1}] = 0$ . But  $\max [[G']_I^s] > 0$  so  $\max [[G']_I^s] > \max [[G']_I^{s+i-1}]$ .
- To show  $\langle \text{rel}(\theta(A_i)), \max [[G']_I^i] \rangle \preceq \langle \text{rel}(A_i), \max [[G]_I^i] \rangle$  for all  $i \in [1, s-1]$ . Let  $i \in [1, s-1]$ .
  - Suppose  $[[G']_I^i] \neq \emptyset$ . Then there exists (5) a grounding substitution  $\sigma$  for  $G'$  and (6) a covering  $D \in \text{covers}_{|,G'}(\theta(A_i))$  such that (7)  $I \models \sigma(D)$  and (8)  $|\sigma(\theta(A_i))| + 1 = \max [[G']_I^i]$ . By (6) it follows that there exists a  $C \subseteq \{A_1, \dots, A_{i-1}, A_{i+1}, \dots, A_m\}$  such that  $C \in \text{covers}_{|,G}(A_i)$  and either:
    - \*  $A_s \in C$  and  $D = \theta(C) \setminus \{\theta(A_s)\} \cup \theta(\{B_1, \dots, B_n\})$  or
    - \*  $A_s \notin C$  and  $D = \theta(C)$ .
 In the first case, by (7)  $I \models \sigma(\theta(C) \setminus \{\theta(A_s)\} \cup \theta(\{B_1, \dots, B_n\}))$  and consequently, since  $I$  is a model for  $c$ ,  $I \models \sigma(\theta(C) \setminus \{\theta(A_s)\} \cup \{\theta(H)\})$ . But then  $I \models \sigma(\theta(C))$  because  $\theta(H) = \theta(A_s)$ . In the second case, by (7),  $I \models \sigma(\theta(C))$ .  
 By (1)  $\sigma \circ \theta$  is a grounding substitution for  $A_1, \dots, A_{s-1}, A_{s+1}, \dots, A_m$ . Let  $\kappa$  be a grounding substitution for  $\sigma(\theta(A_s))$ . Then  $\kappa \circ \sigma \circ \theta$  is a grounding substitution for  $G$  such that  $I \models \kappa(\sigma(\theta(C)))$  since  $\kappa(\sigma(\theta(C))) = \sigma(\theta(C))$ . Hence  $|\kappa(\sigma(\theta(A_i)))| + 1 \in [[G]_I^i]$  and as a result  $\max [[G']_I^i] \leq \max [[G]_I^i]$ , hence  $[[G']_I^i]$  is finite.
  - Suppose  $[[G']_I^i] = \emptyset$ . Then  $[[G']_I^i]$  is trivially finite. Moreover  $\max [[G']_I^i] = 0 \leq \max [[G]_I^i]$ .
- To show  $\langle \text{rel}(\theta(A_i)), \max [[G']_I^{i+n-1}] \rangle \preceq \langle \text{rel}(A_i), \max [[G]_I^i] \rangle$  for all  $i \in [s+1, m]$ . This is analogous to the previous case.

Thus  $[[G']_I]$  is well-defined and  $[[G']_I] \prec_{mul} [[G]_I]$  as required.  $\square$

The proof of theorem 4.2 then follows as immediately consequence of lemma 4.1 since  $\prec_{mul}$  is well-ordered. This completes the termination argument for semi-delay recurrent programs.

## 5. The transformation: termination, soundness and completeness

The force of theorem 4.2 is that it provides a way of developing a program transformation for deriving a correct and reasonably efficient program from a logical specification. The idea is to transform a given program into one which is semi-delay recurrent, but with equivalent declarative semantics. Then, by adding safe delay declarations, a program is obtained which terminates for all goals using semi-local selection.

The transformation is defined at two levels: a transformation that separately acts on each clause in the original program, to yield a clause that is semi-delay recurrent; and a transformation that is applied to each predicate in the original program, to produce an additional new clause. The clause-level transformation replaces each clause defining a predicate  $p \in \Pi$ , with another defining a predicate  $p^{sdr} \in \Pi^{sdr}$ . The predicate-level transformation, when acting on a predicate  $p$ , introduces a clause that defines  $p$  in terms of  $p^{sdr}$  and an auxiliary predicate  $p^{depth} \in \Pi^{depth}$ . The symbol sets  $\Pi$ ,  $\Pi^{sdr}$  and  $\Pi^{depth}$  are required to be disjoint, that is,  $\Pi \cap \Pi^{sdr} = \Pi \cap \Pi^{depth} = \Pi^{sdr} \cap \Pi^{depth} = \emptyset$ . In order to define the transformation, it is convenient to represent depth counters (numbers) as successor terms, that is, encode 0, 1, 2 and arbitrary  $k$  by the terms 0, succ(0), succ(succ(0)) and  $\text{succ}^k(0)$ . Using this representation, the clause-level semi-delay recurrent (SDR) transformation is defined as follows:

**Definition 5.1.** If  $c = H :- B_1, \dots, B_n \in P$  then  $sdr(c) = p^{sdr}(\vec{t}, d) :- B', B'_1, \dots, B'_n$  where

$$B' = \begin{cases} \text{true} & \text{if } \bigwedge_{i=1}^n H \not\neq B_i \\ d = \text{succ}(d') & \text{otherwise} \end{cases} \quad B'_i = \begin{cases} B_i & \text{if } H \not\neq B_i \\ p_i^{sdr}(\vec{t}_i, d') & \text{otherwise} \end{cases}$$

$H = p(\vec{t})$ ,  $B_i = p_i(\vec{t}_i)$  and the variables  $d, d'$  are fresh, that is,  $\{d, d'\} \cap \text{vars}(c) = \emptyset$ .

The predicate true always succeeds; it is a no-operation that is used to simplify the transformation. The unification  $d = \text{succ}(d')$  decrements  $d$  to obtain  $d'$ . It succeeds only if  $d$  equals or exceeds  $\text{succ}(0)$ . For purposes of efficiency, this operation is actually realised as two separate operations: a test and a decrement operating on conventional integers. The predicate-level transformation is defined by:

**Definition 5.2.** If  $p \in \Pi$  then  $sdr(p) = p(x_1, \dots, x_k) :- p^{depth}(x_1, \dots, x_k, d), p^{sdr}(x_1, \dots, x_k, d)$  where the variables  $x_1, \dots, x_k, d$  are distinct and  $k$  is the arity of  $p$ .

Together, the  $sdr$  transformation defines the set of clauses  $P^{sdr} = \{sdr(c) \mid c \in P\} \cup \{sdr(p) \mid p \in \Pi\}$ . The transformation is completed by augmenting  $P^{sdr}$  with a set of clauses  $P^{depth}$  that specify the predicates within  $\Pi^{depth}$ . The intuition is that each predicate  $p^{depth}(x_1, \dots, x_k, d)$  calculates a bound  $d$ , from the arguments  $x_1, \dots, x_k$ , that controls the depth to which the predicate  $p^{sdr}$  is searched. For the program  $P^{sdr} \cup P^{depth}$  to yield finite derivations under semi-local selection, the clauses of  $\Pi^{depth}$  are required to be semi-delay recurrent. The termination issue is considered further in section 5.1. For the program  $P^{sdr} \cup P^{depth}$  to be correct relative to  $P$ , it must preserve each answer that  $P$  can produce (it must be complete relative to  $P$ ) whilst not introducing additional answers (it must be sound relative to  $P$ ). These correctness issues are addressed in section 5.2; in particular the section specifies criteria on the predicates of  $P^{depth}$  that ensure completeness (soundness follows almost immediately). The remainder of this section illustrates the transformation.

**Example 5.1.** Consider again the quicksort program of figure 1, denoting the  $i^{\text{th}}$  clauses of predicates qsort and append by  $\text{qsort}_i$  and  $\text{append}_i$  respectively. Then

$$\begin{aligned} sdr(\text{qsort}_1) &= \text{qsort}^{sdr}([], [], D) :- \text{true}. \\ sdr(\text{qsort}_2) &= \text{qsort}^{sdr}([X|Xs], Ys, D) :- D = \text{succ}(D1), \\ &\quad \text{part}(Xs, X, L, B), \text{qsort}^{sdr}(L, Ls, D1), \text{qsort}^{sdr}(B, Bs, D1), \text{append}(Ls, [X|Bs], Ys) \\ sdr(\text{append}_1) &= \text{append}^{sdr}([], X, X, D) :- \text{true}. \\ sdr(\text{append}_2) &= \text{append}^{sdr}([X|Xs], Ys, [X|Zs], D) :- D = \text{succ}(D1), \text{append}^{sdr}(Xs, Ys, Zs, D1). \end{aligned}$$

and  $sdr(qsort) = qsort(X, Y) :- qsort^{depth}(X, Y, D), qsort^{sdr}(X, Y, D)$   
 $sdr(append) = append(X, Y, Z) :- append^{depth}(X, Y, Z, D), append^{sdr}(X, Y, Z, D)$

### 5.1. The termination issue

The following theorem addresses the termination issue, by showing that the semi-delay recurrency of  $P^{depth}$  is sufficient to ensure that  $P^{sdr} \cup P^{depth}$  is also semi-delay recurrent. The value of this result is that theorem 4.2 guarantees that the program  $P^{sdr} \cup P^{depth}$  will yield finite derivations under semi-local selection, providing that it is equipped with delay declarations that are safe.

**Theorem 5.1.** Let  $I$  be a model of  $P^{sdr} \cup P^{depth}$  and suppose  $P^{depth}$  is semi-delay recurrent wrt a level mapping  $|\cdot|_1$  and  $I$ . Then  $P^{sdr} \cup P^{depth}$  is semi-delay recurrent wrt  $|\cdot|_2$  and  $I$  where

$$|p^{sdr}(\vec{t}, t')|_2 = |t'|_{term-size} \quad |p(\vec{t})|_2 = 0 \quad |p^{depth}(\vec{t}, t')|_2 = |p^{depth}(\vec{t}, t')|_1$$

Thus the program  $P^{sdr} \cup P^{depth}$  will be terminating for all goals under semi-local selection if, for each predicate, a delay declaration is added which is safe wrt  $|\cdot|_2$ . In fact, since the goals  $p(\vec{t})$  are always bounded wrt  $|\cdot|_2$ , delay declarations need only be supplied for the  $p^{sdr}$  and  $p^{depth}$  predicates.

#### Proof:

Let  $\simeq$  be the mutual dependence relation induced from  $P^{sdr} \cup P^{depth}$ . Since  $P^{depth}$  is semi-delay recurrent wrt  $|\cdot|_1$  and  $I$ , it remains to consider each  $c \in P^{sdr}$ :

- Suppose  $c = p(\vec{x}) :- p^{depth}(\vec{x}, d), p^{sdr}(\vec{x}, d)$ . Then  $c$  is semi-delay recurrent wrt  $|\cdot|_2$  and  $I$  since  $p \not\simeq p^{depth}$  and  $p \not\simeq p^{sdr}$ .
- Suppose  $c = p^{sdr}(\vec{t}, d) :- \text{true}, B_1, \dots, B_n$  and  $p \not\simeq rel(B_1), \dots, p \not\simeq rel(B_n)$ . Then again  $c$  is trivially semi-delay recurrent wrt  $|\cdot|_2$  and  $I$ .
- Suppose  $c = p^{sdr}(\vec{t}, d) :- d = \text{succ}(d'), B'_1, \dots, B'_n$ . Let  $i \in [1, n]$  such that  $B'_i = p_i^{sdr}(\vec{t}, d')$ , that is,  $p \simeq p_i$ . Then  $mdcovers_{|\cdot|_2, c}(B'_i) = \{d = \text{succ}(d')\}$  and  $mdcovers_{|\cdot|_2, c}(d = \text{succ}(d')) = \{\emptyset\}$ , hence  $covers_{|\cdot|_2, c}(B'_i) = \{d = \text{succ}(d')\}$ . Let  $\theta$  be a grounding substitution for  $c$  such that  $I \models \theta(d = \text{succ}(d'))$ . Hence  $\theta(d) = \theta(\text{succ}(d'))$ , thus  $|\theta(p^{sdr}(\vec{t}, d))|_2 = |\theta(d)|_{term-size} = |\theta(\text{succ}(d'))|_{term-size} = 1 + |\theta(d')|_{term-size} > |\theta(d')|_{term-size} = |\theta(p_i^{sdr}(\vec{t}, d'))|_2$ .

□

The proof of theorem 5.1 exploits the property that  $d = \text{succ}(d')$  is the only cover of each recursive body atom  $B'_i = p_i^{sdr}(\vec{t}, d')$  occurring within the body of a clause defining  $p^{sdr}$ . Thus, once  $d = \text{succ}(d')$  is resolved, then the remaining goals  $B'_1, \dots, B'_n$  can be scheduled in any order. Thus coroutining can be liberally applied.

The main consequence of theorem 5.1 is that a program can be transformed into another that is semi-delay recurrent and whose semantics are equivalent in a sense which will be examined shortly.

## 5.2. The issues of soundness and completeness

The transformation is completed by providing clauses  $P^{depth}$  for the predicates  $\Pi^{depth}$ . These clauses compute depth bounds on the subsequent computation, for example, the rôle of the  $qsort^{depth}$  predicate is to establish a bound on the depth of the recursive calls to  $qsort^{sdr}$ . The bound needs to be large enough to enable the search space to be probed to sufficient depth to find all the answers to a query. For example, the goal  $qsort([1, 2, 3], L)$  will fail if  $qsort^{depth}([1, 2, 3], L, D)$  instantiates  $D$  to  $succ(0)$ , since more than one recursive call to the predicate  $qsort^{sdr}$  is required to find the unique answer  $L = [1, 2, 3]$ . On the other hand, the goal  $qsort([1, 2, 3], L)$  will not produce an incorrect answer. Thus the depth bound effects the completeness of the transformation rather than soundness. Nevertheless, to ensure that completeness is preserved, the definition of  $qsort^{depth}$  needs to be suitably constrained. The constraint manifests itself as a condition for completeness in the equivalence theorem that is stated below.

**Theorem 5.2.** Let  $I$  and  $J$  be minimal models for  $P$  and  $P^{sdr} \cup P^{depth}$ .

- Then if  $J \models p(\vec{t})$  and  $p \in \Pi$  it follows that  $I \models p(\vec{t})$  (soundness of the transformation).
- Suppose that whenever  $J \models p^{sdr}(\vec{t}, succ^j(0))$  and  $J \not\models p^{sdr}(\vec{t}, succ^{j-1}(0))$  there exists  $k \geq 0$  such that  $J \models p^{depth}(\vec{t}, succ^{j+k}(0))$ . Then if  $I \models p(\vec{t})$  it follows that  $J \models p(\vec{t})$  (completeness of the transformation).

The criteria for completeness is a statement about the relative sizes of the depth parameters of  $p^{sdr}$  and  $p^{depth}$ . The condition  $J \not\models p^{sdr}(\vec{t}, succ^{j-1}(0))$  ensures that  $j$  is minimal, that is, no smaller  $j$  exists such that  $J \models p^{sdr}(\vec{t}, succ^j(0))$ . The theorem does not require these depth parameters to coincide; instead it is sufficient for the depth parameter of  $p^{depth}$  to match or exceed that of  $p^{sdr}$ . This is a useful result from a practical point of view since upper bounds are easier to compute than exact bounds.

In order to reason about the equivalence of the minimal models  $I$  and  $J$ , and in particular prove theorem 5.2, the standard  $T_P$  [23] is introduced since its iterative formulation provides a way of constructing an inductive equivalence argument.

**Definition 5.3.**

$$T_P(I) = \left\{ \theta(H) \left| \begin{array}{l} c = H:-B_1, \dots, B_n \in P \quad \wedge \\ \theta \text{ is a grounding substitution for } c \quad \wedge \\ \theta(B_1), \dots, \theta(B_n) \in I \end{array} \right. \right\}$$

The soundness argument is, in fact, more involved than the completeness argument. Since  $P^{sdr} \cup P^{depth}$  defines intermediate predicates  $\Pi^{sdr}$  that do not appear in  $P$ , if  $T_P^k(\emptyset) \models p(\vec{t})$  it does not necessarily follow that  $T_{P^{sdr} \cup P^{depth}}^k(\emptyset) \models p(\vec{t})$ ; in general, a larger value  $k' > k$  is required for  $T_{P^{sdr} \cup P^{depth}}^{k'}(\emptyset) \models p(\vec{t})$ . To construct an inductive hypothesis for soundness, and specifically define  $k'$  in terms of  $k$ , it is necessary to reason about the relative positions (or *level*) of predicates with the call graph hierarchy of  $P$ . This requirement leads to the following definition.

**Definition 5.4.** Suppose that  $\simeq$  and  $\sqsupseteq$  are the relations over  $\Pi = \{p_1, \dots, p_n\}$  induced from  $P$ . Then define  $level(p_i) = \min\{l_i \mid \langle l_1, \dots, l_n \rangle \in level(\Pi)\}$  where

$$level(\Pi) = \left\{ \langle l_1, \dots, l_n \rangle \in \mathbb{N}^n \left| \begin{array}{l} p_j \simeq p_k \rightarrow l_j = l_k \quad \wedge \\ p_j \not\simeq p_k \wedge p_j \sqsupseteq p_k \rightarrow l_j > l_k \end{array} \right. \right\}$$

**Example 5.2.** Returning to the quicksort program of figure 1 and example 4.6,  $\Pi = \{\text{qsort}, \text{part}, \text{append}, \text{leq}, \text{gt}\}$  and using the notation in definition 5.4,

$$\text{level}(\Pi) = \{ \langle l_1, \dots, l_5 \rangle \in \mathbb{N}^5 \mid l_1 > l_2 \wedge l_1 > l_3 \wedge l_1 > l_4 \wedge l_1 > l_5 \wedge l_2 > l_4 \wedge l_2 > l_5 \}$$

hence  $\text{level}(\text{leq}) = 1$ ,  $\text{level}(\text{gt}) = 1$ ,  $\text{level}(\text{append}) = 1$ ,  $\text{level}(\text{part}) = 2$  and  $\text{level}(\text{qsort}) = 3$ .

The proof for theorem 5.2 (which follows immediately after that of lemma 5.2) is constructed in three stages. Firstly, lemma 5.1 states a monotonicity condition on  $p^{sdr}$  (its proof is straightforward and thus is not included). Secondly, lemma 5.1 is then used with the *level* machinery to argue a form of equivalence between  $T_P$  and  $T_{P^{sdr} \cup P^{depth}}$ . This gives lemma 5.2. Thirdly, the equivalence result of lemma 5.2 is then used to establish theorem 5.2 in conjunction with the property that the minimal model of a program  $P$  coincides with the limit of  $T_P$ . Interestingly, this final stage of proof requires that the clauses defined in  $P^{depth}$  to not invoke any predicates defined within  $P^{sdr}$ , which is a reasonable requirement. This construction completes the argument that the declarative semantics of the original and transformed programs (restricted to the predicates defined in the original program) coincide.

**Lemma 5.1.** If  $T_{P^{sdr} \cup P^{depth}}^k(\emptyset) \models p^{sdr}(\vec{t}, d)$  then  $T_{P^{sdr} \cup P^{depth}}^k(\emptyset) \models p^{sdr}(\vec{t}, \text{succ}(d))$  for all  $p \in \Pi$  and for all  $k \geq 0$ .

**Lemma 5.2.** Let  $J$  and  $M$  be minimal models for  $P^{sdr} \cup P^{depth}$  and  $P^{depth}$ . Suppose that whenever  $J \models p^{sdr}(\vec{t}, \text{succ}^j(0))$  and  $J \not\models p^{sdr}(\vec{t}, \text{succ}^{j-1}(0))$  there exists  $l \geq 0$  such that  $J \models p^{depth}(\vec{t}, \text{succ}^{j+l}(0))$ .

- Then if  $T_P^k(\emptyset) \models p(\vec{t})$  it follows that

- $T_{P^{sdr}}^{k+\text{level}(p)}(M) \models p^{sdr}(\vec{t}, \text{succ}^l(0))$  for some  $l \geq 0$ ;
- $T_{P^{sdr}}^{k+\text{level}(p)+1}(M) \models p(\vec{t})$

- Then if  $T_{P^{sdr}}^k(M) \models p^{sdr}(\vec{t}, \text{succ}^l(0)), p(\vec{t})$  for some  $l \geq 0$  it follows that  $T_P^k(\emptyset) \models p(\vec{t})$ .

**Proof:**

Consider the first case, as the second is straightforward. Proof by induction. Suppose  $T_P^k(\emptyset) \models p(\vec{t})$ .

- To show that  $T_{P^{sdr}}^{k+\text{level}(p)}(M) \models p^{sdr}(\vec{t}, \text{succ}^l(0))$  for some  $l \geq 0$ . Since  $T_P^k(\emptyset) = \emptyset$  the result vacuously holds for  $k = 0$ . Thus consider  $k \geq 1$ . Suppose  $T_P^k(\emptyset) \models p(\vec{t})$ . Then a grounding substitution  $\theta$  exists for a clause  $c = H :- B_1, \dots, B_n \in P$  such that  $\theta(H) = p(\vec{t})$  and  $\theta(B_1), \dots, \theta(B_n) \in T_P^{k-1}(\emptyset)$ . Suppose  $H = p(\vec{s})$  and  $B_i = p_i(\vec{s}_i)$ .
  - Suppose  $sdr(c) = p^{sdr}(\vec{s}, d) :- d = \text{succ}(d'), B'_1, \dots, B'_n$ . Let  $I = \{i \in [1, n] \mid p_i \simeq p\}$ .
    - \* Consider  $i \in I$ . By induction it follows that  $T_{P^{sdr}}^{(k-1)+\text{level}(p)}(M) \models p_i^{sdr}(\theta(\vec{s}_i), d_i)$  for some  $d_i = \text{succ}^{l_i}(0)$  since  $\text{level}(p_i) = \text{level}(p)$ . Since  $I \neq \emptyset$ , put  $l = \max(\{l_i \mid i \in I\})$  and define  $\sigma = \theta \cup \{d = \text{succ}(\text{succ}^l(0)), d' \mapsto \text{succ}^l(0)\}$ . By lemma 5.1 it follows that  $T_{P^{sdr}}^{(k-1)+\text{level}(p)}(M) \models \sigma(p_i^{sdr}(\vec{s}_i, d'))$ .
    - \* Consider  $i \in [1, n] \setminus I$ . By induction it follows that  $T_{P^{sdr}}^{(k-1)+\text{level}(p_i)+1}(M) \models p_i(\theta(\vec{s}_i))$ . But  $\text{level}(p_i) < \text{level}(p)$  hence  $T_{P^{sdr}}^{(k-1)+\text{level}(p)}(M) \models p_i(\sigma(\vec{s}_i))$ .



Because  $\theta$  is a grounding substitution for  $c$ , it follows that  $\sigma$  is a grounding substitution for  $sdr(c)$ . Hence  $T_{P^{sdr}}^{k+level(p)}(M) \models \sigma(p^{sdr}(\vec{s}, d)) = p^{sdr}(\vec{t}, succ^{l+1}(0))$  as required.

- Suppose  $sdr(c) = p^{sdr}(\vec{s}, d) : -\text{true}, B'_1, \dots, B'_n$ . The argument is a simplified version of the previous case.
- To show  $T_{P^{sdr}}^{k+level(p)+1}(M) \models p(\vec{t})$ . Recall that  $c = p(\vec{x}) : -p^{depth}(\vec{x}, d), p^{sdr}(\vec{x}, d) \in P^{sdr}$ . By the preceding argument there exists a minimal  $l \geq 0$  such that  $T_{P^{sdr}}^{k+level(p)}(M) \models p^{sdr}(\vec{t}, succ^l(0))$ . Now there exists  $k \geq 0$  such that  $M \models p^{depth}(\vec{t}, succ^{k+l}(0))$ . By lemma 5.1,  $T_{P^{sdr}}^{k+level(p)}(M) \models p^{sdr}(\vec{t}, succ^{k+l}(0))$ . A grounding substitution  $\theta$  for  $c$  exists such that  $\theta(\vec{x}) = \vec{t}$  and  $\theta(d) = succ^{k+l}(0)$ . Hence  $T_{P^{sdr}}^{k+level(p)+1}(M) \models p(\vec{t})$  as required.

□

**Proof:**

Let  $I, J$  and  $M$  be minimal models for  $P, P^{sdr} \cup P^{depth}$  and  $P^{depth}$ .

- Suppose  $J \models p(\vec{t})$ . There exists  $k \geq 0$  such that  $T_{P^{sdr} \cup P^{depth}}^k(\emptyset) \models p(\vec{t})$  and because  $P^{depth}$  is not defined using  $P^{sdr}$ , it follows that  $T_{P^{sdr}}^k(M) \models p(\vec{t})$ . Since  $T_{P^{sdr}}^k(M) \models p(\vec{t})$  it follows  $T_{P^{sdr}}^k(M) \models p^{sdr}(\vec{t}, succ^l(0))$  for some  $l \geq 0$ . Whence  $T_P^k(\emptyset) \models p(\vec{t})$  by lemma 5.2. But  $I \supseteq T_P^k(\emptyset)$  hence the result follows.
- Suppose  $I \models p(\vec{t})$ . There exists  $k \geq 0$  such that  $T_P^k(\emptyset) \models p(\vec{t})$ , whence  $T_{P^{sdr}}^{k+level(p)+1}(M) \models p(\vec{t})$  by lemma 5.2. But  $J \supseteq T_{P^{sdr} \cup P^{depth}}^{k+level(p)+1}(\emptyset) = T_{P^{sdr}}^{k+level(p)+1}(M)$  again since the predicates in  $P^{depth}$  are not defined in terms of  $P^{sdr}$ . The result then follows.

□

The astute reader will notice that although the transformed program is semi-delay recurrent, there is no guarantee that an arbitrary query to the program will not suspend. However, it has recently been shown [17] that, given a logic program with delay declarations, it is possible to infer classes of goal that do not lead to suspensions. The reader is referred to [17, section 8.3] which reports (among other things) that the transformed version of mergesort will not suspend for  $\text{msort}(L, S)$  when  $L$  or  $S$  are ground.

## 6. The transformation: automation and efficiency

The development of the transformation, thus far, has focussed on the issues of termination, soundness and completeness. This section shows how program analysis, when augmented with program transformation, can be used to infer bounds on the depths of search trees. This relieves the programmer from much of the burden of applying the transformation. This section also examines the efficiency of the resulting code.

### 6.1. Automating the discovery of depth bounds

The construction of  $P^{sdr}$  is straightforward to automate. Automatically synthesising definitions for  $P^{depth}$  is more challenging: theorem 5.2 asserts that although soundness of the transformation follows

immediately, for completeness the  $P^{depth}$  predicates need to be defined so as to provide a bound on the depth of a search tree. Computing such a bound cannot be achieved in general since it is well-known from computability theory that there are partial recursive functions which cannot be extended to total recursive functions. The problem of computing a bound is analogous to the problem of checking termination of a given program and goal; although the problem is unsolvable in the general case, many specific cases are susceptible to automation [14].

One well-established technique in termination checking is argument size analysis [4, 15] that infers size relationships between the arguments of predicates. In this context, argument size relationships are used to observe a decrease between the size of arguments in the head and the size of the arguments in the body. This analysis cannot be applied directly to  $P^{sdr}$  to infer a relationship of the form  $d \leq \sum_{i=1}^k c_i |t_i|$  for each predicate  $p^{sdr}(t_1, \dots, t_k, d)$  where  $|\cdot|$  is a norm and  $c_1, \dots, c_k \in \mathbb{R}$  is a set of coefficients. In fact no such invariant holds for  $P^{sdr}$  no matter how  $P^{depth}$  is defined. By lemma 5.1 it follows that if  $J$  is the minimal model of  $P^{sdr} \cup P^{depth}$  and  $J \models p(\vec{t}, d)$  then  $J \models p(\vec{t}, \text{succ}(d))$ . The subtlety is that it is not the syntactic form of the  $d \leq \sum_{i=1}^k c_i |t_i|$  inequality that is inappropriate for ensuring completeness, it is its interpretation. To satisfy the completeness criteria of theorem 5.2 it is sufficient to infer an invariant  $d \leq \sum_{i=1}^k c_i |t_i|$  that holds in the sense that if  $J \models p^{sdr}(\vec{t}, \text{succ}^d(0))$  and  $J \not\models p^{sdr}(\vec{t}, \text{succ}^{d-1}(0))$  then  $d \leq \sum_{i=1}^k c_i |t_i|$ . Fortunately such invariants can be inferred automatically by applying argument size analysis to a transformation of  $P$ , named  $P^{min}$ , that is engineered solely for the purpose of deducing these invariants. The presentation of this approach is structured in four phases: section 6.1.1 defines the transformation that yields  $P^{min}$ ; section 6.1.2 explains the precise relationship between the programs  $P^{depth}$  and  $P^{min}$ ; section 6.1.3 shows how this relationship can be exploited to infer bounds on the depth of search trees; and section 6.1.4 explains how  $P^{depth}$  can be constructed from depth bound invariants.

### 6.1.1. A transformation for discovering depth bounds

As with  $P^{sdr}$ , the transformation that constructs  $P^{min}$  is defined at two levels: a transformation that separately acts on each clause in the original program and a transformation that is applied to each predicate in the original program, to produce an additional new clause. The clause-level transformation replaces each clause defining a predicate  $p \in \Pi$ , with another defining a predicate  $p^{min} \in \Pi^{min}$  where  $\Pi \cap \Pi^{min} = \emptyset$ . This clause-level and predicate-level transformations are defined as follows:

**Definition 6.1.** If  $c = H :- B_1, \dots, B_n \in P$  then  $min(c) = p^{min}(\vec{t}, d) :- B', B'_1, \dots, B'_n$  where

$$B' = \begin{cases} d = 0 & \text{if } \bigwedge_{i=1}^n H \not\approx B_i \\ \max(d_1, \dots, d_n, d'), d = \text{succ}(d') & \text{otherwise} \end{cases} \quad B'_i = \begin{cases} B_i, d_i = 0 & \text{if } H \not\approx B_i \\ p_i^{min}(\vec{t}_i, d_i) & \text{otherwise} \end{cases}$$

$H = p(\vec{t})$ ,  $B_i = p_i(\vec{t}_i)$ , and the variables  $d, d', d_i$  are fresh, that is,  $\{d, d', d_1, \dots, d_n\} \cap \text{vars}(c) = \emptyset$ .

**Definition 6.2.** If  $p \in \Pi$  then  $min(p) = p(x_1, \dots, x_k) :- p^{min}(x_1, \dots, x_k, d)$  where the variables  $x_1, \dots, x_k, d$  are distinct and  $k$  is the arity of  $p$ .

For any given clause  $c \in P$ , the clause  $min(c)$  can be considered to be an instrumented version of  $sdr(c)$  that counts the minimal number of times that a loop has to be traversed to yield an answer. For any non-recursive clause the loop count is zero. For any recursive clause, the loop count is exactly one plus the maxima of the loop count of each of the recursive body atoms. Thus the transformation is formulated

```

max(x, x) :- true.
max(0, x2, x) :- max(x2, x).
max(x1, 0, x) :- max(x1, x).
max(succ(x1), succ(x2), succ(x)) :- max(x1, x2, x).
max(0, x2, x3, x) :- max(x2, x3, x).
max(x1, 0, x3, x) :- max(x1, x3, x).
max(x1, x2, 0, x) :- max(x1, x2, x).
max(succ(x1), succ(x2), succ(x3), succ(x)) :- max(x1, x2, x3, x).

```

Figure 6. The max predicates for arity 2, 3 and 4 cases; other cases are defined analogously

```

d(N, _, 0) :- nat(N).
d(X, X, 1).
d(F*N, X, N*DF) :- nat(N), d(F, X, DF).
d(N*F, X, N*DF) :- nat(N), d(F, X, DF).
d(-F, X, -DF) :- d(F, X, DF).
d(X#succ(N), X, succ(N)*X#N) :- nat(N).
d(sin(X), X, cos(X)).
d(cos(X), X, -sin(X)).
d(F*G, X, DF*G+DG*F) :- d(F, X, DF), d(G, X, DG).
d(F/G, X, (DF*G-DG*F)/(G*G)) :- d(F, X, DF), d(G, X, DG).
d(F+G, X, DF+DG) :- d(F, X, DF), d(G, X, DG).
d(F-G, X, DF-DG) :- d(F, X, DF), d(G, X, DG).
nat(0).
nat(succ(N)) :- nat(N).

```

Figure 7. Symbolic differentiation in Prolog

in terms of a series of auxiliary predicates  $\max(t_1, t)$ ,  $\max(t_1, t_2, t)$ ,  $\dots$ ,  $\max(t_1, t_2, \dots, t_k, t)$ , etc that are defined to hold iff  $t_i = \text{succ}^{n_i}(0)$ ,  $t = \text{succ}^n(0)$  and  $n = \max\{n_1, \dots, n_k\}$ . Definitions for these predicated are listed in figure 6. These definitions constitute a set of clauses that henceforth will be denoted as  $P^{max}$ .  $P^{max}$  is assumed to include max predicates of arity  $2, \dots, k+1$  where  $k$  is the maximum number of body atoms that occur in any clause of  $P$ . Moreover, the max predicate symbol is required to be unique, that is,  $\max \notin \Pi$  and  $\max \notin \{p^{min} \mid p \in \Pi\}$ . The transformation  $min(p)$  that operates on a predicate  $p$  merely introduces a clause that discards the count. Together, these two transformation define the set of clauses  $P^{min} = \{min(c) \mid c \in P\} \cup \{min(p) \mid p \in \Pi\}$ .

**Example 6.1.** Suppose that  $P$  is the program listed in figure 7.  $P$  defines a predicate  $d(t_1, t_2, t_3)$  that differentiates a function  $t_1$  wrt  $t_2$  to obtain the derivative  $t_3$ . The rules (clauses) in the program correspond to rules of symbolic differentiation, for example, the ninth and tenth clauses of  $P$  respectively implement the product and quotient rules. The operator  $\#$  represents power, thus a query such as  $d(x\#\text{succ}(\text{succ}(0)), x, F)$  instantiates  $F$  to  $\text{succ}(\text{succ}(0)) * x\#\text{succ}(0)$ . The predicate  $\text{nat}(t)$  holds iff

- (1)  $d(X1, X2, X3) :- d^{min}(X1, X2, X3, \_).$
- (2)  $d^{min}(N, \_, 0, D) :- D = 0, \text{nat}(N).$
- (3)  $d^{min}(X, X, 1, D) :- D = 0.$
- (4)  $d^{min}(F * N, X, N * DF, D) :- D = \text{succ}(D1), \text{nat}(N), d^{min}(F, X, DF, D1).$
- (5)  $d^{min}(N * F, X, N * DF, D) :- D = \text{succ}(D1), \text{nat}(N), d^{min}(F, X, DF, D1).$
- (6)  $d^{min}(-F, X, -DF, D) :- D = \text{succ}(D1), d^{min}(F, X, DF, D1).$
- (7)  $d^{min}(X \# \text{succ}(N), X, \text{succ}(N) * X \# N, D) :- D = 0, \text{nat}(N).$
- (8)  $d^{min}(\sin(X), X, \cos(X), D) :- D = 0.$
- (9)  $d^{min}(\cos(X), X, -\sin(X), D) :- D = 0.$
- (10)  $d^{min}(F * G, X, DF * G + DG * F, D) :-$   
 $D = \text{succ}(D1), \max(D2, D3, D1), d^{min}(F, X, DF, D2), d^{min}(G, X, DG, D3).$
- (11)  $d^{min}(F / G, X, (DF * G - DG * F) / (G * G), D) :-$   
 $D = \text{succ}(D1), \max(D2, D3, D1), d^{min}(F, X, DF, D2), d^{min}(G, X, DG, D3).$
- (12)  $d^{min}(F + G, X, DF + DG, D) :-$   
 $D = \text{succ}(D1), \max(D2, D3, D1), d^{min}(F, X, DF, D2), d^{min}(G, X, DG, D3).$
- (13)  $d^{min}(F - G, X, DF - DG, D) :-$   
 $D = \text{succ}(D1), \max(D2, D3, D1), d^{min}(F, X, DF, D2), d^{min}(G, X, DG, D3).$
- (14)  $\text{nat}(X) :- \text{nat}^{min}(X1, \_).$
- (15)  $\text{nat}^{min}(0, D) :- D = 0.$
- (16)  $\text{nat}^{min}(\text{succ}(N), D) :- D = \text{succ}(D1), \text{nat}^{min}(N, D1).$
- (17)  $\max(X, X).$
- (18)  $\max(0, X2, X) :- \max(X2, X).$
- (19)  $\max(X1, 0, X3) :- \max(X1, X).$
- (20)  $\max(\text{succ}(X1), \text{succ}(X2), \text{succ}(X)) :- \max(X1, X2, X).$

Figure 8.  $P^{min} \cup P^{max}$  for the symbolic differentiation program  $P$ 

$t$  is a natural number represented in successor notation; it is used to detect constants. The program  $P^{min} \cup P^{max}$  is listed in figure 8. For purposes of presentation, this program is actually a slight specialisation of  $P^{min} \cup P^{max}$ . For instance, if  $c = d(F * N, X, N * DF) :- \text{nat}(N), d(F, X, DF)$  then

$$\min(c) = \begin{cases} d^{min}(F * N, X, N * DF, D) :- \\ \max(D1, D2, D3), D = \text{succ}(D3), \text{nat}(N), D1 = 0, d^{min}(F, X, DF, D2) \end{cases}$$

Since  $\max(0, D2, D3)$  holds iff  $\max(D2, D3)$  holds iff  $D2 = D3$  holds, then  $\min(c)$  can be specialised to the form given in figure 8. The desire for brevity explains the cascaded structure of the max predicates.

### 6.1.2. The relationship between the depth bounds in $P^{sdr}$ and $P^{min}$

This section provides a bridging result that reinterprets the completeness criteria stated in theorem 5.2 in terms of requirements on the program  $P^{min} \cup P^{max}$ . Before the key result, corollary 6.1, is given and shown, a supporting lemma is stated. The lemma relates the program  $P^{min} \cup P^{max}$  to the program  $P^{sdr} \cup P^{depth}$  subject to the condition that the predicates in  $P^{depth}$  are vacuous, that is, they do not

impose any depth constraints. The lemma formalises the intuition that the counters in the program  $P^{min}$  compute the absolute minimal number of times that a loop has to be traversed to yield an answer (the proof of the lemma is straightforward and thus is not included).

**Lemma 6.1.** Suppose  $J$  and  $K$  are minimal models for  $P^{sdr} \cup P^{depth}$  and  $P^{min} \cup P^{max}$  respectively where  $P^{depth} = \{p^{depth}(x_1, \dots, x_k, d) \mid p \in \Pi\}$ . Let  $p \in \Pi$  then

- $K \models p(\vec{t})$  iff  $J \models p(\vec{t})$ ;
- $K \models p^{min}(\vec{t}, d)$  iff  $d = \text{succ}^{\min(L)}(0)$  and  $L = \{l \geq 0 \mid J \models p^{sdr}(\vec{t}, \text{succ}^l(0))\} \neq \emptyset$ .

**Corollary 6.1.** Let  $I$  and  $K$  be minimal models for  $P$  and  $P^{min} \cup P^{max}$  respectively. Suppose that for each  $p \in \Pi$  there exist functions  $f_1, \dots, f_m \in U_P^n \rightarrow \mathbb{N}$  such that if  $K \models p^{min}(\vec{t}, \text{succ}^k(0))$  then

$$k \leq f_1(\vec{t}), \dots, k \leq f_m(\vec{t})$$

Suppose  $P^{depth}$  is defined so that if  $J$  is the minimal model of  $P^{sdr} \cup P^{depth}$  then if  $J \models p^{depth}(\vec{t}, d)$  then  $d = \text{succ}^k(0)$  and  $\min\{f_1(\vec{t}), \dots, f_m(\vec{t})\} \leq k$ . Then it follows that if  $I \models p(\vec{t})$  then  $J \models p(\vec{t})$ .

**Proof:**

Suppose  $J \models p^{sdr}(\vec{t}, \text{succ}^j(0))$  and  $J \not\models p^{sdr}(\vec{t}, \text{succ}^{j-1}(0))$ . By lemma 6.1,  $K \models p^{min}(\vec{t}, \text{succ}^j(0))$ . But if  $J \models p^{depth}(\vec{t}, d)$  then  $d = \text{succ}^l(0)$  and  $j \leq \min\{f_1(\vec{t}), \dots, f_m(\vec{t})\} \leq l$ , hence there exists  $k = l - j \geq 0$  such that  $J \models p^{depth}(\vec{t}, \text{succ}^{j+k}(0))$  and the result follows by theorem 5.2.  $\square$

The significance of corollary 6.1 is that it provides a condition for completeness that is formulated in terms of  $P^{min} \cup P^{max}$ . In particular, there is no reason why  $f_i(\vec{t})$  cannot take the form  $\lfloor b + \sum_{i=1}^n c_i |t_i| \rfloor$  where  $b, c_i \in \mathbb{Q}$  and  $|\cdot|$  is a norm. Bounds of this form can be extracted from  $P^{min} \cup P^{max}$  by argument-size analysis. The corollary also provides a specification for  $P^{depth}$ ; for completeness it is sufficient to realise  $p^{depth}$  so that if  $p^{depth}(\vec{t}, \text{succ}^k(0))$  holds then  $k$  is clamped below by  $\min\{f_1(\vec{t}), \dots, f_m(\vec{t})\}$ .

### 6.1.3. Extracting depth bounds from $P^{min}$

The corollary reduces the problem of inferring depth bounds for  $P^{sdr}$  to the problem of deducing argument-size relationships on the minimal model of  $P^{min} \cup P^{max}$ . The established approach to finding such invariants [4, 15] involves characterising the minimal model of a  $\text{CLP}(\mathcal{R})$  program. The  $\text{CLP}(\mathcal{R})$  program is derived from the  $P^{min} \cup P^{max}$  program in such a way that the minimal model of the former, say  $L$ , describes the minimal model of the latter, say  $K$ . This notion of description is a formal concept:  $L$  describes  $K$  with respect to some norm  $|\cdot|$  if whenever  $K \models p(\vec{t})$  it follows that  $L \models p(|\vec{t}|)$  where  $|\langle t_1, \dots, t_n \rangle| = \langle |t_1|, \dots, |t_n| \rangle$ . The value of this concept is that it ensures that any invariant on  $L$  can be safely reinterpreted as a size invariant on  $K$ . For example, if  $\vec{c} \cdot \vec{n} \leq b$  for all  $L \models p(\vec{n})$  then it follows that  $\vec{c} \cdot |\vec{t}| \leq b$  for all  $K \models p(\vec{t})$ . This is convenient because  $L$  is a simpler computational domain than  $K$ . If  $L \models p(\vec{n})$  then  $\vec{n}$  is merely a vector of numbers, however if  $K \models p(\vec{t})$  then  $\vec{t}$  is a vector of terms of arbitrary depth.

The  $\text{CLP}(\mathcal{R})$  program has to be derived so as to ensure that  $L$  describes  $K$ . This derivation proceeds by replacing each syntactic equation  $t_1 = t_2$  in the  $P^{min} \cup P^{max}$  program with a linear equation  $\vec{c} \cdot \vec{x} = b$  in the  $\text{CLP}(\mathcal{R})$  program. The equation  $\vec{c} \cdot \vec{x} = b$  describes  $t_1 = t_2$  wrt  $|\cdot|$  iff whenever  $\theta$  is a grounding substitution for  $t_1 = t_2$  such that  $\theta(t_1) = \theta(t_2)$  then  $\vec{c} \cdot |\theta(\vec{x})| = b$ .

- (1)  $d(A, B, C) :- d^{min}(A, B, C, D), 1 \leq A, 1 \leq B, 1 \leq C, 1 \leq D.$
- (2)  $d^{min}(A, B, C, D) :- C=1, D=1, \text{nat}(A), 1 \leq A, 1 \leq B, 1 \leq D, 1 \leq C.$
- (3)  $d^{min}(A, B, C, D) :- A=B, C=1, D=1, 1 \leq A, 1 \leq D, 1 \leq B, 1 \leq C.$
- (4)  $d^{min}(A, B, C, D) :- A=1+E+F, C=1+F+G, D=1+H, \text{nat}(F), d^{min}(E, B, G, H),$   
 $1 \leq E, 1 \leq F, 1 \leq B, 1 \leq G, 1 \leq D, 1 \leq H, 1 \leq A, 1 \leq C.$
- (5)  $d^{min}(A, B, C, D) :- A=1+E+F, C=1+E+G, D=1+H, \text{nat}(E), d^{min}(F, B, G, H),$   
 $1 \leq E, 1 \leq F, 1 \leq B, 1 \leq G, 1 \leq D, 1 \leq H, 1 \leq A, 1 \leq C.$
- (6)  $d^{min}(A, B, C, D) :- A=1+E, C=1+F, D=1+G, d^{min}(E, B, F, G), 1 \leq E, \dots, 1 \leq C.$
- (7)  $d^{min}(A, B, C, D) :- A=2+B+E, C=3+B+2*E, D=1, \text{nat}(E), 1 \leq B, \dots, 1 \leq C.$
- (8)  $d^{min}(A, B, C, D) :- A=1+B, C=1+B, D=1, 1 \leq B, 1 \leq D, 1 \leq A, 1 \leq C.$
- (9)  $d^{min}(A, B, C, D) :- A=1+B, C=2+B, D=1, 1 \leq B, 1 \leq D, 1 \leq A, 1 \leq C.$
- (10)  $d^{min}(A, B, C, D) :- A=1+E+F, C=3+E+F+G+H, D=1+I, \max(J, K, I),$   
 $d^{min}(E, B, G, J), d^{min}(F, B, H, K), 1 \leq E, \dots, 1 \leq C.$
- (11)  $d^{min}(A, B, C, D) :- A=1+E+F, C=5+E+3*F+G+H, D=1+I, \max(J, K, I),$   
 $d^{min}(E, B, G, J), d^{min}(F, B, H, K), 1 \leq E, \dots, 1 \leq C.$
- (12)  $d^{min}(A, B, C, D) :- A=1+E+F, C=1+G+H, D=1+I, \max(J, K, I),$   
 $d^{min}(E, B, G, J), d^{min}(F, B, H, K), 1 \leq E, \dots, 1 \leq C.$
- (14)  $\text{nat}(A) :- \text{nat}^{min}(B, C), 1 \leq A, 1 \leq B, 1 \leq C.$
- (15)  $\text{nat}^{min}(A, B) :- A=1, B=1, 1 \leq B, 1 \leq A.$
- (16)  $\text{nat}^{min}(A, B) :- A=1+C, B=1+D, \text{nat}^{min}(C, D), 1 \leq C, \dots, 1 \leq A.$
- (17)  $\max(A, B) :- A=B, 1 \leq A, 1 \leq B.$
- (18)  $\max(A, B, C) :- A=1, \max(B, C), 1 \leq B, 1 \leq C, 1 \leq A.$
- (19)  $\max(A, B, C) :- B=1, \max(A, C), 1 \leq A, 1 \leq C, 1 \leq B.$
- (20)  $\max(A, B, C) :- A=1+D, B=1+E, C=1+F, \max(D, E, F), 1 \leq D, \dots, 1 \leq C.$

Figure 9. Term-size description of  $P^{min} \cup P^{max}$ 

**Example 6.2.** The linear equation  $C = 3 + X + 2 * N$  describes  $C = \text{succ}(N) * X \# N$  wrt the term size norm  $|\cdot|_{\text{term-size}}$ . To see this, let  $\theta$  be a grounding substitution of the syntactic equation. Then  $\theta(C) = \text{succ}(\theta(N)) * \theta(X) \# \theta(N)$ , hence  $|\theta(C)|_{\text{term-size}} = |\text{succ}(\theta(N)) * \theta(X) \# \theta(N)|_{\text{term-size}} = 3 + |\theta(X)|_{\text{term-size}} + 2|\theta(N)|_{\text{term-size}}$ . The linear equation expresses the relative sizes of the any ground instance of the variables  $C, X$  and  $N$  that satisfies the syntactic equation.

The concept of description can be lifted from equations to clauses by augmenting the concept with normalisation, that is, the process of replacing of terms that arise in head and body atom arguments with fresh variables. Thus a clause  $c$  describes  $p(\vec{t}): - p_1(\vec{t}_1) \dots p_n(\vec{t}_n)$  wrt a norm  $|\cdot|$  iff  $c$  takes the form  $p(\vec{x}): - e, e_1, \dots, e_n, p_1(\vec{t}_1) \dots p_n(\vec{t}_n)$  where  $\vec{x}$  and  $\vec{x}_i$  are vectors of fresh variables and the linear equations  $e$  and  $e_i$  describe the syntactic equations  $\vec{x} = \vec{t}$  and  $\vec{x}_i = \vec{t}_i$  wrt the norm  $|\cdot|$ . A program  $P_1$  describes another  $P_2$  wrt a norm  $|\cdot|$  iff for each clause  $c_2 \in P_2$  there exists a clause  $c_1 \in P_1$  such that  $c_1$  describes  $c_2$  wrt the norm  $|\cdot|$ .

**Example 6.3.** The program listed in figure 9 describes the program  $P^{min} \cup P^{max}$  wrt the term-size norm

$|\cdot|_{\text{term-size}}$ . Notice that each clause  $c$  is augmented with inequalities  $1 \leq y_1, \dots, 1 \leq y_m$  that assert the positivity of the variables  $y_1, \dots, y_m$  that occur in the clause. These inequalities do not alter the meaning of the program because if  $\theta$  is a grounding substitution for  $c$  then  $1 \leq |\theta(y_i)|_{\text{term-size}}$  for each variable  $y_i$ . The inequalities are traditionally included to make the program less sensitive to the precision loss that can arise in fixpoint computation. Note, too, that clause 12 in figure 9 describes clause 12 and 13 in figure 8. The clauses differ only in the names of functor symbols and term-size description does not preserve information about the particular functors that occur in syntactic equations.

Let  $Eqn_X$  denote the set of linear equalities and inequalities defined over a finite set of variables  $X$ . Each element of  $Eqn_X$  takes the form of either  $\vec{c}_1 \cdot \vec{y} = b_1$  or  $\vec{c}_2 \cdot \vec{y} \leq b_2$  where  $|\vec{c}_i| = |\vec{y}|$ ,  $b_i \in \mathbb{Z}$  and the elements of  $\vec{c}_i$  and  $\vec{y}$  are drawn from  $\mathbb{Z}$  and  $X$  respectively. The set of integer solutions for  $\vec{c}_1 \cdot \vec{y} = b_1$ ,  $\vec{c}_2 \cdot \vec{y} \leq b_2$  and any linear system  $E \subseteq Eqn_X$  are formally defined as follows:

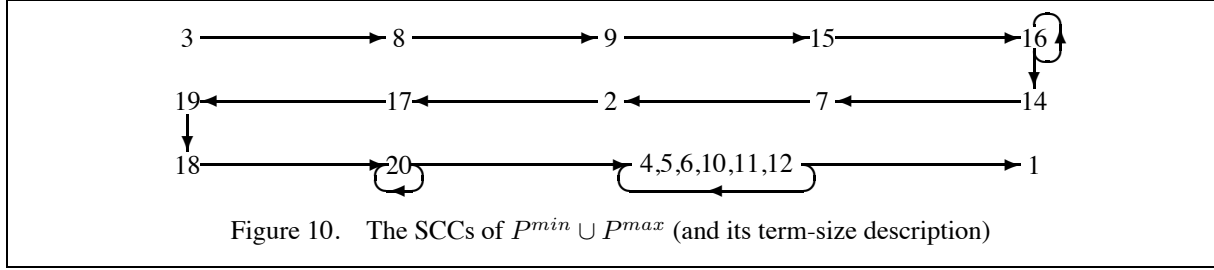
**Definition 6.3.** Suppose  $\vec{x} = \langle x_1, \dots, x_n \rangle$  and  $\vec{y} = \langle y_1, \dots, y_m \rangle$ . Then

$$\text{soln}_{\vec{x}}(\vec{c} \cdot \vec{y} \leq b) = \left\{ \langle r_1, \dots, r_n \rangle \in \mathbb{Z}^n \mid \begin{array}{l} \vec{c} \cdot \langle r'_1, \dots, r'_m \rangle \leq b \quad \wedge \\ r_i = r'_j \text{ for all } x_i = y_j \end{array} \right\}$$

The set  $\text{soln}_{\vec{x}}(\vec{c} \cdot \vec{y} = b)$  is defined analogously and  $\text{soln}_{\vec{x}}(E) = \bigcap_{e \in E} \text{soln}_{\vec{x}}(e)$ .

Two linear systems  $E_1, E_2 \subseteq Eqn_X$  are partially ordered by the subset relation on their solution sets, that is,  $E_1 \models E_2$  iff  $\text{soln}_{\vec{x}}(E_1) \subseteq \text{soln}_{\vec{x}}(E_2)$  where  $\text{var}(\vec{x}) = X$ . One subtlety of using (finite) linear systems as a computational domain is that infinite ascending chains can arise. Consider, for example, the sequence of linear systems  $E_1, E_2, \dots$  where  $E_i = \{\langle 1 \rangle \cdot \langle y \rangle = 0, \langle -1 \rangle \cdot \langle x \rangle \leq 0, \langle 1 \rangle \cdot \langle x \rangle \leq i\}$ . Each  $E_i$  is a finite set that defines a line  $i$  units long that is rooted at the origin and runs parallel to the  $x$  axis. This sequence forms a chain since  $E_i \models E_{i+1}$  and  $E_{i+1} \not\models E_i$  (each  $E_{i+1}$  is strictly larger than its predecessor). A monotonic operator over a domain that admits only finite chains is guaranteed to reach a fixpoint in a finite number of steps and thus terminate. Termination of such an operator, however, may be compromised if the domain contains infinite chains. In the case of  $Eqn_X$ , the simplest way to enforce converge is to work within a sub-domain that contains only finite chains. One such sub-domain is constructed from finite sets of inequalities drawn from  $\text{Mono}_X = \{x \leq 0, 0 \leq x, x \leq y, x + y \leq z, x \leq y + z \mid x, y, z \in X\}$ .  $\text{Mono}_X$  is so named because it resembles the monotonic domain that is used for observing termination in Datalog programs [7]. If  $E_1, E_2 \subseteq Eqn_X$  then  $E_1$  and  $E_2$  can *both* be described using  $\text{Mono}_Y$  by applying the operation  $E_1 \vee_Y E_2$  that is defined by  $E_1 \vee_Y E_2 = \{e \in \text{Mono}_Y \mid E_1 \models \{e\} \wedge E_2 \models \{e\}\}$ . (This operation that conservatively describes two abstractions with another is called merge [18]). The system  $E_1 \vee_Y E_2$  is finite if  $Y$  is finite.

To capture linear invariants between the arguments of predicates, it is necessary to lift the  $\models$  ordering to atoms paired with linear systems as follows  $\langle p(\vec{x}_1), E_1 \rangle \models \langle p(\vec{x}_2), E_2 \rangle$  iff  $\text{soln}_{\vec{x}_1}(E_1) \subseteq \text{soln}_{\vec{x}_2}(E_2)$ . Observe that two pairs  $\langle p(\vec{x}_1), E_1 \rangle$  and  $\langle p(\vec{x}_2), E_2 \rangle$  that differ syntactically may express the same invariants, that is,  $\langle p(\vec{x}_1), E_1 \rangle \models \langle p(\vec{x}_2), E_2 \rangle$  and  $\langle p(\vec{x}_2), E_2 \rangle \models \langle p(\vec{x}_1), E_1 \rangle$ . To express invariants between argument positions it is thus necessary to construct sets of syntactically different but equivalence pairs. (This is more than an aesthetic predilection since this construction simplifies the way formal arguments are matched against actual arguments.) Formally, equivalence is defined by  $\langle p(\vec{x}_1), E_1 \rangle \equiv \langle p(\vec{x}_2), E_2 \rangle$  iff  $\langle p(\vec{x}_1), E_1 \rangle \models \langle p(\vec{x}_2), E_2 \rangle$  and  $\langle p(\vec{x}_2), E_2 \rangle \models \langle p(\vec{x}_1), E_1 \rangle$  which, in turn, induces a notion of equivalence class. To simultaneously record the invariants that hold on different predicates, the ordering is



further extended to sets of equivalence classes that contain exactly one class  $[\langle p(\vec{x}), E \rangle]_{\equiv}$  for each  $p \in \Pi$  by defining  $I_1 \models I_2$  iff  $\langle p(\vec{x}), E_1 \rangle \models \langle p(\vec{x}), E_2 \rangle$  whenever  $[\langle p(\vec{x}), E_1 \rangle]_{\equiv} \in I_1$  and  $[\langle p(\vec{x}), E_2 \rangle]_{\equiv} \in I_2$ . Infinite chains can be avoided under this ordering by ensuring that the sets only contain equivalence classes of the form  $[\langle p(\vec{x}), E \rangle]_{\equiv}$  where  $E \subseteq Mono_{var(\vec{x})}$ . Under this assumption, the sets of equivalence classes constitute a finite lattice; finiteness follows because  $\Pi$  is finite and each  $p \in \Pi$  has finite arity. The bottom element of this structure is given by  $\perp = \{[\langle p(\vec{x}), false \rangle]_{\equiv} \mid p \in \Pi\}$  where *false* denotes any unsatisfiable system, that is,  $soln_{\vec{x}}(false) = \emptyset$ . Sets of equivalence classes provide a computation domain for the following operator that is designed to discover invariants in  $CLP(\mathcal{R})$  programs:

**Definition 6.4.**

$$T_c^{CLP}(I) = \left\{ [\langle p(\vec{x}), F \vee_{var(\vec{x})} F' \rangle]_{\equiv} \mid \begin{array}{l} c = p(\vec{x}) :- E, p_1(\vec{x}_1), \dots, p_n(\vec{x}_n) \quad \wedge \\ [\langle p_i(\vec{x}_i), E_i \rangle]_{\equiv} \in I \wedge [\langle p(\vec{x}), F \rangle]_{\equiv} \in I \wedge \\ F' = E \cup (\cup_{i=1}^n E_i) \end{array} \right\}$$

The operator is defined for individual  $CLP(\mathcal{R})$  clauses of the form  $p(\vec{x}) :- E, p_1(\vec{x}_1), \dots, p_n(\vec{x}_n)$  where  $E$  is any finite linear system. The  $\vee_{var(\vec{x})}$  operation ensures that any computed equivalence class has a representative member  $[\langle p(\vec{x}), E \rangle]$  satisfying the requirement  $E \subseteq Mono_{var(\vec{x})}$ . This operator can be lifted to the level of a program  $P = \{c_1, \dots, c_n\}$  by defining  $T_P^{CLP}(I) = I_{n+1}$  where  $I_1 = I$  and  $I_{i+1} = T_{c_i}^{CLP}(I_i)$ . Since  $T_P^{CLP}$  is monotonic it follows that the least fixpoint  $\text{lfp}(T_P^{CLP})$  exists and since  $T_P^{CLP}$  is continuous it follows that  $\text{lfp}(T_P^{CLP})$  can be computed as the limit of the sequence  $I_{k+1} = T_P^{CLP}(I_k)$  where  $I_1 = \perp$ . The following proposition explains the way in which the least fixpoint faithfully describes the minimal model of the original program. The proof is not given since it can be constructed straightforwardly by adapting proofs that have been reported elsewhere [3, 18].

**Proposition 6.1.** Suppose a  $CLP(\mathcal{R})$  program  $P_1$  describes a program  $P_2$  wrt a norm  $|\cdot|$ . If  $K$  is the minimal model of  $P_2$  and  $K \models p(\vec{t})$  then  $[\langle p(\vec{x}), E \rangle]_{\equiv} \in \text{lfp}(T_{P_1}^{CLP})$  and  $|\vec{t}| \in soln_{\vec{x}}(E)$ .

The way  $T_P^{CLP}$  in which is formulated in terms of  $T_{c_i}^{CLP}$  provides a route towards efficient fixpoint computation. Observe that if a set of clauses  $P = \{c_1, \dots, c_n\}$  is ordered so that each  $c_1, \dots, c_{k-1}$  do not depend on  $c_k$ , that is,  $c_1 \not\supseteq c_k, \dots, c_{k-1} \not\supseteq c_k$ , then the fixpoint calculation reduces to  $I_{k+1} = T_{c_k}^{CLP}(I_k)$  where  $I_1 = \perp$ . More generally,  $P$  can be partitioned into a set of SCCs  $\{S_1, \dots, S_m\}$ . If the set is ordered so that  $c_{k-1} \not\supseteq c_k$  for all  $c_k \in S_k$  and for all  $c_{k-1} \in \cup_{i=1}^{k-1} S_i$  then  $\text{lfp}(T_P^{CLP}) = J_{m+1}$  where  $J_{k+1}$  is a limit calculated for each individual  $S_k$  and  $J_1 = \perp$ . The limit  $J_{k+1}$  is, in turn, defined by  $I_{i+1} = T_{S_k}^{CLP}(I_i)$  where  $I_1 = J_k$ . Moreover, if  $S_k = \{c_k\}$  and  $c_k \not\supseteq c_k$  then  $J_{k+1} = T_{c_k}^{CLP}(J_k)$ .



$k$	$p(\vec{x})$	$F'$ (truncated)	$F' \vee_{var(\vec{x})} F'$
3	$d^{min}(A, B, C, D)$	$A \leq B+D, B \leq A, C \leq A, A \leq B, D \leq C, C \leq D$	$A \leq B+D, B \leq A, C \leq A, A \leq B, D \leq C, C \leq D$
8	$d^{min}(A, B, C, D)$	$A \leq C+D, B+D \leq A, A \leq B+D, B+D \leq C, C \leq B+D, D \leq B$	$A \leq B+D, B \leq A, C \leq A, D \leq B, D \leq C$
9	$d^{min}(A, B, C, D)$	$A \leq C+D, B+D \leq A, A \leq B+D, A+D \leq C, C \leq A+D, C \leq A+B$	$A \leq B+D, C \leq A+D, B \leq A, D \leq B, D \leq C$
15	$nat^{min}(A, B)$	$0 \leq A, 0 \leq B, A \leq B, B \leq A$	$0 \leq A, 0 \leq B, A \leq B, B \leq A$
16	$nat^{min}(A, B)$	$0 \leq A, 0 \leq B, A \leq B, B \leq A$	$0 \leq A, 0 \leq B, A \leq B, B \leq A$
16	$nat^{min}(A, B)$	$0 \leq A, 0 \leq B, A \leq B, B \leq A$	$0 \leq A, 0 \leq B, A \leq B, B \leq A$
14	$nat(A)$	$0 \leq A$	$0 \leq A$
7	$d^{min}(A, B, C, D)$	$B+D \leq A, A+D \leq C, D \leq B, 0 \leq D$	$A \leq B+C, B \leq A, D \leq B, D \leq C, 0 \leq D$
2	$d^{min}(A, B, C, D)$	$C \leq A, C \leq B, D \leq C, C \leq D, 0 \leq D$	$D \leq A, D \leq B, D \leq C, 0 \leq D$
17	$max(A, B)$	$0 \leq A, 0 \leq B, A \leq B, B \leq A$	$0 \leq A, 0 \leq B, A \leq B, B \leq A$
19	$max(A, B, C)$	$A \leq B+C, B \leq A, C \leq A, A \leq C$	$A \leq B+C, B \leq A, C \leq A, A \leq C$
18	$max(A, B, C)$	$B \leq A+C, A \leq B, C \leq B, B \leq C$	$C \leq A+B, A \leq C, B \leq C$
20	$max(A, B, C)$	$C \leq A+B, A \leq C, B \leq C$	$C \leq A+B, A \leq C, B \leq C$
20	$max(A, B, C)$	$C \leq A+B, A \leq C, B \leq C$	$C \leq A+B, A \leq C, B \leq C$
4	$d^{min}(A, B, C, D)$	$D \leq A, D \leq C, 0 \leq D, 0 \leq B$	$D \leq A, D \leq C, 0 \leq D, 0 \leq B$
5	$d^{min}(A, B, C, D)$	$D \leq A, D \leq C, 0 \leq D, 0 \leq B$	$D \leq A, D \leq C, 0 \leq D, 0 \leq B$
6	$d^{min}(A, B, C, D)$	$D \leq A, D \leq C, 0 \leq D, 0 \leq B$	$D \leq A, D \leq C, 0 \leq D, 0 \leq B$
10	$d^{min}(A, B, C, D)$	$A+D \leq C, D \leq A, 0 \leq D, 0 \leq B$	$D \leq A, D \leq C, 0 \leq D, 0 \leq B$
11	$d^{min}(A, B, C, D)$	$A+D \leq C, D \leq A, 0 \leq D, 0 \leq B$	$D \leq A, D \leq C, 0 \leq D, 0 \leq B$
12	$d^{min}(A, B, C, D)$	$D \leq A, D \leq C, 0 \leq D, 0 \leq B$	$D \leq A, D \leq C, 0 \leq D, 0 \leq B$
4	$d^{min}(A, B, C, D)$	$D \leq A, D \leq C, 0 \leq D, 0 \leq B$	$D \leq A, D \leq C, 0 \leq D, 0 \leq B$
1	$d(A, B, C)$	$0 \leq A, 0 \leq B, 0 \leq C$	$0 \leq A, 0 \leq B, 0 \leq C$

Figure 11. Fixpoint computation for the term-size description of  $P^{min} \cup P^{max}$ 

**Example 6.4.** Figure 10 presents the 14 SCCs of the program  $P^{min} \cup P^{max}$  (and its description). The SCCs are  $\{3\}$ ,  $\{8\}$ ,  $\dots$ ,  $\{4,5,6,10,11,12\}$ ,  $\{1\}$ . There are just 3 non-trivial SCCs: one the contains clause 16, another that contains clause 20 and another that contains clauses 4, 5, 6, 10, 11 and 12. The forward arrows in figure 10 denote a total ordering on the  $S_k$  such that  $c_{k-1} \not\sqsupseteq c_k$  for all  $c_k \in S_k$  and for all  $c_{k-1} \in \cup_{i=1}^{k-1} S_i$ . The back arrows indicate where iteration is required in the fixpoint calculation.

**Example 6.5.** Figure 11 presents the fixpoint calculation for the description of  $P^{min} \cup P^{max}$  given in figure 9. Iteration commences with  $I_1 = \perp$ . The first column specifies the clause number  $k$  used to calculate  $I_{i+1} = T_{c_k}^{CLP}(I_i)$ . The second and third columns give the atom  $p(\vec{x})$  and the linear system that is  $F'$  obtained by evaluating the clause  $c_k$  with  $I_i$ . For brevity only those inequalities and equalities of  $F'$  that relate to  $\vec{x}$  are listed (the truncated version of  $F'$  formally corresponds to the projection of  $F'$  onto  $\vec{x}$ ). The third column gives  $F' \vee_{var(\vec{x})} F'$ . The horizontal bars delineate the start and finish of the iterative sub-calculations for the 3 SCCs. Note that to check stability for the clauses  $\{4,5,6,10,11,12\}$ , it is not necessary reevaluate the entire SCC. If the clauses are visited in the same order, it is sufficient to observe that a clause does not alter the iterate on reevaluation. Observe the  $D \leq A$  and  $D \leq C$  invariants inferred for  $d^{min}$  and that  $B \leq A$  is inferred for  $nat^{min}$ .

**Example 6.6.** Whether an invariant can be discovered depends on its particular syntactic form. This is illustrated by the quicksort and mergesort programs listed in figures 1 and 2. It is only when the inequalities of  $Mono_X$  are enriched to include constants of 1, 0 or -1 that the following can be inferred:

$$\begin{aligned}
&\langle \text{qsort}^{min}(A, B, C), \{1 \leq C, A \leq B \leq A, C \leq A, C \leq B\} \rangle \\
&\langle \text{part}^{min}(A, B, C, D, E), \{1 + A \leq C + D, E \leq A, 1 \leq D, 1 \leq C, 1 \leq E, 1 \leq B\} \rangle \\
&\langle \text{append}^{min}(A, B, C, D), \{D \leq A, D \leq C, 1 + C \leq A + B \leq 1 + C, 1 \leq B, 1 \leq D\} \rangle \\
\\
&\langle \text{msort}^{min}(A, B, C), \{1 \leq C, A \leq B, B \leq A, C \leq A\} \rangle \\
&\langle \text{merge}^{min}(A, B, C, D), \{1 \leq A, 1 \leq B, 1 \leq D, D \leq C, 1 + C \leq A + B \leq 1 + C\} \rangle \\
&\langle \text{split}^{min}(A, B, C, D), \{1 \leq C, 1 \leq D, D \leq B, D \leq 1 + C, B + C \leq 1 + A \leq B + C\} \rangle
\end{aligned}$$

Richer classes of invariant could be inferred by augmenting or replacing  $Mono_X$  with systems of inequalities with arbitrary coefficients that involve two variables [38], or arbitrary affine equations [15] or arbitrary systems of linear inequations [4]. In fact there is no reason why the depth bounds have to be linear and another route towards automation is offered by recent work on inferring non-linear loop invariants [36]. In this context, an invariant is a conjunction of polynomial equalities: the values that variables can take at a program point is described by the roots of polynomials over those variables. For example, the abstraction  $\langle x + y, y(y + 1) \rangle$  represents those sets of points in 2-dimensional space that satisfy both  $x + y = 0$  and  $y(y + 1) = 0$ , namely, the points  $\langle 0, 0 \rangle$  and  $\langle 1, -1 \rangle$ . The elegance of the scheme is that merge coincides with the intersection of polynomial ideals. Since loops can entail an unbounded number of merges, the ideals occurring within loops can contain polynomials of arbitrary degree. A (widening) operator is thus introduced to enforce converge. An implementation of the method suggests that interesting invariants can be calculated in an automatic way for non-trivial programs. Moreover, there seems no reason why this method cannot be adapted to infer non-linear argument size relationships.

The invariants given above for the differentiation, quicksort and mergesort programs were derived automatically. The main components of this system are a polyhedral library (800 LOC) that implements the operations on linear systems; an abstracter unit (800 LOC) that computes a description of a program wrt either  $|\cdot|_{term-size}$  or  $|\cdot|_{list-length}$ ; and a meta-interpreter (150 LOC) that calculates the fixpoint. The polyhedral library uses the CLP(Q) library of SICStus 3 to decide  $E \models \{\vec{c} \cdot \vec{x} \leq b\}$  where  $E$  is a linear system. This decision procedure is required to detect stability and also compute  $E_1 \vee_Y E_2$ . The CLP(Q) library is useful since  $E \models \{\vec{c} \cdot \vec{x} \leq b\}$  holds iff  $E \cup \{(-\vec{c}) \cdot \vec{x} \leq -(b + 1)\}$  is unsatisfiable (whenever  $\vec{c}$ ,  $\vec{x}$  and  $b$  are integral). For further implementation details, the reader is referred to the sources that can be downloaded from the homepage of the first author.

#### 6.1.4. Synthesising $P^{depth}$ from depth bounds

Not all the invariants inferred by the  $T_P^{CLP}$  operator yield depth bounds: some may not involve the depth bound and others may not bound it from above. To illustrate exactly how  $P^{depth}$  can be constructed from the inferred invariants suppose  $[\langle p^{min}(\vec{x}), E \rangle]_{\equiv} \in \text{lfpt}(T_P^{CLP})$  where  $P$  describes  $P^{min} \cup P^{max}$  wrt the term size norm  $|\cdot|_{term-size}$  and consider  $\vec{c} \cdot \vec{y} \leq b \in E$ . Without loss of generality it is sufficient to consider  $\vec{c}' \cdot \vec{x} \leq b$  where  $\text{soln}_{\vec{x}}(\vec{c}' \cdot \vec{x} \leq b) = \text{soln}_{\vec{x}}(\vec{c} \cdot \vec{y} \leq b)$ . The invariant  $\vec{c}' \cdot \vec{x} \leq b$  provides an upper bound on the depth iff  $0 < c'_n$  where  $n$  is the arity of  $p^{min}$ . It then follows that  $x_n \leq (b - \sum_{i=1}^{n-1} c'_i x_i) / c'_n$ . By proposition 6.1 it follows that if  $K$  is the minimal model of  $P^{min} \cup P^{max}$  and  $K \models p(\vec{t})$  then

```

d_depth(F, X, R, D) :- term_size(F, D1, D), term_size(R, D2, D), d_aux(D1, D2, D).
:- block d_aux(-, -, ?).
d_aux(D1, D2, D) :- nonvar(D1) -> D is D1 - 1; D is D2 - 1.
nat_depth(N, D) :- term_size(N, D, D).
term_size(Term, Size, D) :- term_size(Term, 0, Size, D).
:- block term_size(-, ?, ?, -), term_size(?, -, ?, -).
term_size(Term, Acc, Size, D) :-
    nonvar(D) -> true; Term =.. [_ | Args], Acc1 is Acc + 1, term_sizes(Args, Acc1, Size, D).
term_sizes([], Size, Size, _).
term_sizes([Arg | Args], Acc, Size, D) :-
    term_size(Arg, Acc, NewAcc, D), term_sizes(Args, NewAcc, Size, D).

```

Figure 12. Term size predicates

$|t_n|_{term-size} \leq (b - \sum_{i \neq 1}^{n-1} c'_i |x_i|_{term-size}) / c'_n$ . By the construction of  $P^{min} \cup P^{max}$ ,  $t_k = \text{succ}^d(0)$ , hence  $|\text{succ}^d(0)|_{term-size} = 1 + d$  and therefore  $d \leq f(t_1, \dots, t_{n-1}) = \lfloor (b - \sum_{i \neq 1}^{n-1} c'_i |t_i|_{term-size}) / c'_n \rfloor - 1$  which gives a function  $f$  compliant with the completeness requirements of corollary 6.1.

**Example 6.7.** Returning to the differentiation program, from figure 11 consider the pairs  $\langle \text{nat}^{min}(\vec{x}), E_1 \rangle$  and  $\langle d^{min}(\vec{y}), E_2 \rangle$  where  $\vec{x} = \langle A, B \rangle$ ,  $\vec{y} = \langle A, B, C, D \rangle$  and

$$E_1 = \left\{ \begin{array}{ll} \langle -1, 0 \rangle \cdot \vec{x} \leq 0, & \langle 0, -1 \rangle \cdot \vec{x} \leq 0, \\ \langle 1, -1 \rangle \cdot \vec{x} \leq 0, & \langle -1, 1 \rangle \cdot \vec{x} \leq 0 \end{array} \right\} \quad E_2 = \left\{ \begin{array}{ll} \langle -1, 0, 0, 1 \rangle \cdot \vec{y} \leq 0, & \langle 0, 0, -1, 1 \rangle \cdot \vec{y} \leq 0, \\ \langle 0, 0, 0, -1 \rangle \cdot \vec{y} \leq 0, & \langle 0, -1, 0, 0 \rangle \cdot \vec{y} \leq 0 \end{array} \right\}$$

For  $E_1$  the last inequality gives  $f(t) = \lfloor (0 - (-1|t|_{term-size})) / 1 \rfloor - 1 = |t|_{term-size} - 1$  and for  $E_2$  the first and second inequalities yield  $f_1(t_1, t_2, t_3) = |t_1|_{term-size} - 1$  and  $f_2(t_1, t_2, t_3) = |t_3|_{term-size} - 1$ .

The paradox of corollary 6.1 is that it specifies a sufficient condition for completeness on  $P^{depth}$  in terms of its minimal model (which is variable-free) and functions  $f_1, \dots, f_m$  that operate over tuples of terms drawn from the Herbrand universe (which are also variable-free). However, this does not preclude a predicate  $p^{depth}(\vec{t}, d)$  instantiating  $d$  before  $\vec{t}$  is entirely ground. This is because it is sufficient to delay  $p^{depth}$  until  $\min\{f_1(\vec{t}), \dots, f_m(\vec{t})\}$  is bounded on the set of variable-free instances of  $\vec{t}$ . By corollary 6.1 it only necessary to ensure  $\min\{f_1(\vec{t}), \dots, f_m(\vec{t})\} \leq d$ , hence  $d$  can be instantiated to this bound.

**Example 6.8.** In the case of  $d^{min}(t_1, t_2, t_3)$ , the function  $\min\{f_1(t_1, t_2, t_3), f_2(t_1, t_2, t_3)\}$  is bounded on the set of variable-free instances of  $\langle t_1, t_2, t_3 \rangle$  when either  $t_1$  or  $t_3$  are ground. Figure 12 lists code for  $d^{depth}$  based on this tactic. In the case of  $\text{nat}^{depth}$ , although  $f(t) = |t|_{term-size} - 1$ , the definition of the predicate is simplified by assigning  $d = |t|_{term-size}$ . With these predicates thus defined, it is possible to invoke  $d$  in reverse mode to perform symbolic integration. For example, the query  $d(F, x, (-\sin(x)) * \sin(x) - \cos(x) * \cos(x)) / (\sin(x) * \sin(x))$  will instantiate  $F$  to  $\cos(x) / \sin(x)$ .

```

msort(L, S) :- msort_depth(L, S, D), msort_sdr(L, S, D).
msort_depth(L, S, D) :- list_length(L, D1, D), list_length(S, D2, D), msort_aux(D1, D2, D).
:- block msort_aux(-, -, ?).
msort_aux(D1, D2, D) :- D1 = D2, D = D1.
:- block msort_sdr(?, ?, -).
msort_sdr([], [], _).
msort_sdr([X], [X], _).
msort_sdr([X, Y | Xs], S, D) :- 1 ≤ D, D1 is D - 1,
    split_sdr(Xs, L1, L2, D),
    msort_sdr([X|L1], S1, D1), msort_sdr([Y|L2], S2, D1), merge_sdr(S1, S2, S, D).
%:- block split_sdr(?, ?, ?, -).
split_sdr([], [], [], _).
split_sdr([X | L], [X | L1], L2, D) :- 1 ≤ D, D1 is D - 1, split_sdr(L, L2, L1, D1).
%:- block merge_sdr(?, ?, ?, -).
merge_sdr([], Ys, Ys, _).
merge_sdr(Xs, [], Xs, _).
merge_sdr([X|Xs], [Y|Ys], [X|Zs], D) :- 1 ≤ D, D1 is D - 1,
    leq(X, Y), merge_sdr(Xs, [Y|Ys], Zs, D1).
merge_sdr([X|Xs], [Y|Ys], [Y|Zs], D) :- 1 ≤ D, D1 is D - 1,
    gt(X, Y), merge_sdr([X|Xs], Ys, Zs, D1).

```

Figure 13. Optimised version of the control generation transformation for the mergesort program

## 6.2. Efficiency

The essential idea behind the transformation is to ensure termination by delaying possibly non-terminating goals until certain arguments become rigid. This section discusses the efficiency of this approach to control generation. Relative to the original program, the transformation induces additional overheads in terms of synchronisation, depth checking and rigidity detection. We consider each of these overheads in turn. Firstly, suspending and resuming a goal controlled via a block declaration is a constant time operation in SICStus; indeed this operation is directly supported at the level of the abstract machine [10]. Secondly, the test against one and decrement operations both translate to single abstract machine instructions in SICStus 3 [10]. Thirdly, although the transformation introduces rigidity checks that are realised with block declarations, it should be noted that block declarations, or more expensive devices such as a synchronising meta-call [37], are likely to be used in any *ad hoc* attempt at specifying the control.

If efficiency is really a critical issue, then it should be noted that the use of depth bounds – which are merely required to be upper bounds – offers scope for performance tuning. Returning to the mergesort example, recall from section 3 that a depth of  $l$  where  $l$  is the length of the first and second arguments of `msort` is a convenient upper bound on the depth since it avoids the calculation of a transcendental function. However, recall also that `split` is depth bounded by the length of its first argument and `merge` is depth bounded by the sum of the length of its first and second arguments. It therefore follows that  $l$  also bounds the depth of `split` and `merge` as well as that of `msort`. This observation enables rigidity checks

within split and merge to be eliminated and leads to the program given in figure 13. Furthermore, the block declarations for `split_sdr` and `merge_sdr` will always be satisfied and therefore can be removed (for clarity they have been commented out). The optimised version of the program only performs rigidity checking on the initial input; following the call to `msort_aux` the program runs without synchronisation checks and the only overhead is the decrementation and depth testing. The net result is that, with SICStus 3.8.5 on a 1GHz, 256MByte PC, the optimised program averages 75 msec to sort 2048 random numbers. The original program given in figure 2 averages 128 msec on the same data set, whereas the unoptimised transformed version of figure 4 takes 166 msec.

Other sample programs are available from the homepage of the first author. The code `qsort_naive` represents an attempt at making the classic quicksort algorithm reversible. This program averages 618 msec to sort 1024 random numbers but, alas, only produces one solution in reverse mode before looping. The code `qsort_control` generated by a naive application of the transformation averages 838 msec on random data, but is reversible. In reverse mode, it enumerates combinations in a similar fashion to mergesort, when its second argument is instantiated to a list of ordered numbers.

The program `nqueens_control` was obtained by a straightforward application of transformation to the program `nqueens_naive`. The `nqueens_naive` program encodes the  $n$ -queens problem but is actually buggy because its `perm` predicate – which is intended to encode a permutation generator – does not universally terminate for some queries.

```
perm([], []).
perm([X | Xs], P) :- select(P, X, Rs), perm(Xs, Rs).
```

```
select([N | Ns], N, Ns).
select([N | Ns], S, [N | Rs]) :- select(Ns, S, Rs).
```

Specifically, `perm([1, 2, 3], P)` only produces one solution  $P = [1, 2, 3]$  before looping. However, the generated program `nqueens_control` terminates and takes 92 secs to count the number of solutions for the  $n = 10$  board instance. The code is guaranteed to terminate no matter what the ordering of the original body goals, and in particular the `perm(Range, Soln)` and `safe(Soln)` goals in the clause:

```
nqueens(N, Soln):-
    length(Range, N), length(Soln, N), range(Range, 1), perm(Range, Soln), safe(Soln).
```

can be reordered without compromising termination. The goal `safe(Soln)` checks that the board configuration represented by the list `Soln` is safe whereas `range(Range, 1)` instantiates elements of the list `Range` to consecutive numbers starting at 1. Since SICStus implements left-to-right scheduling by default, with the proviso that suspended goals are invoked as soon as possible, this reordering introduces coroutining between the goals `safe(Soln)` and `perm(Range, Soln)`. In effect, bindings to the elements of the list `Soln` are incrementally generated by `perm(Range, Soln)`. As soon as `Soln` contains two elements, representing two queens that can take one another, then `Soln` is an unsafe board configurations; it is unsafe no matter how the elements of `Soln` are instantiated. By ordering `safe(Soln)` before `perm(Range, Soln)`, the goal `safe(Soln)` will check the safety of a new queen, relative to the existing queens, as soon as the queen is added to the board. When the new queen is verified to be safe with respect to the existing queens, control switches back from `safe(Soln)` to `perm(Range, Soln)` to add the next queen. This form of coroutining improves search since unsafe configurations can be detected without necessarily adding all the queens to

the board. In fact, coroutining reduces the running time from 83 secs to 681 msec for  $n = 10$ . It should be noted that transformation does not choose a propitious goal ordering: it merely introduces coroutining capability. In fact, since the transformation operates on pure logic programs, it must ignore any control constructs that already exist in the program. For the purposes of a fair comparison, block declarations were added to `nqueens_naive` to mimic this coroutining behaviour. The termination problem was finessed by observing that `perm(P, [1, 2, 3])` *does* correctly enumerate permutations, hence the arguments in the call `perm(Range, Soln)` were reversed to enforce termination. The resulting program requires 691 msec for  $n = 10$ , which is close to the speed of the transformed code. Finally, notice that the transformation preserves goal ordering, hence if a tester precedes a generator in the original program, then coroutining naturally arises in the transformed program. For example, in the case of a sort defined thus:

```
sort(X,Y) :- ordered(Y), permutation(X,Y).
```

the transformed program will also coroutine, though the additional overheads will incur a 45% slowdown.

Experimental work can never be conclusive, but these results do suggest that the transformation is not prohibitively expensive. Moreover, the resulting code is flexible. Since the transformation is underpinned by semi-delay recurrency (rather than delay recurrency [27]) goals can be reordered to introduce coroutining into the transformed code and thereby improve efficiency. The flexibility offered by the use of upper bounds allows the performance to be tuned by choosing judicious upper bounds that eliminate rigidity checks. However, even without these refinements, the transformation produces code that is surprisingly efficient considering its generality.

## 7. Conclusions

The aim of control generation is to automatically derive a computation rule for a program that is reasonably efficient but does not compromise program correctness. The problem has been effectively tackled by transforming a program into a semantically equivalent one, introducing safe delay declarations and defining a flexible computation rule, which ensures that all goals for the transformed program terminate. Furthermore, it has been shown that the answers computed by the transformed program are complete with respect to the declarative semantics. This is significant.

Beyond the theoretical aspects of the work, its practicality has been demonstrated. In particular, it has been shown how transformed programs can be straightforwardly implemented in a standard logic programming system, and how such a program can be optimised to reduce the number of rigidity checks. Furthermore, with the proposed transformation, the termination problems caused by speculative output bindings are eliminated without the use of a local computation rule or other costly overhead. Moreover, capability for coroutining can contribute to the efficiency of the generated code.

**Acknowledgements** We would like to thank Elena Marchiori for clarifying details of her work and Enric Rodríguez-Carbonell for checking our account of his work. Thanks are also due to Samir Genaim for providing inorder, Axel Simon for help with simplex, and to the referees for their insightful comments.

## References

- [1] Apt, K. R., Pedreschi, D.: Reasoning about Termination of Pure Prolog Programs, *Information and Computation*, **106**(1), 1993, 109–157.
- [2] Apt, K. R., Pedreschi, D.: Modular Termination Proofs for Logic and Pure Prolog programs, *Advances in Logic Programming Theory* (G. Levi, Ed.), Oxford University Press, 1994, Also available as technical report CS-R9316 from Centrum voor Wiskunde en Informatica, CWI, Amsterdam.
- [3] Barbuti, R., Giacobazzi, R., Levi, G.: A General Framework for Semantics-Based Bottom-Up Abstract Interpretation of Logic Programs, *ACM Transactions on Programming Languages and Systems*, **15**(1), 1993, 133–181.
- [4] Benoy, F., King, A.: Inferring Argument Size Relationships with CLP( $\mathcal{R}$ ), *Logic Programming Synthesis and Transformation* (J. P. Gallagher, Ed.), 1207, Springer-Verlag, 1996.
- [5] Bezem, M.: Strong Termination of Logic Programs, *The Journal of Logic Programming*, **15**(1&2), 1993, 79–97.
- [6] Bossi, A., Etalle, S., Rossi, S., J.-G. Smaus: Semantics and Termination of Simply-Moded Logic Programs with Dynamic Scheduling, *ACM Transactions on Computational Logic*, **15**(3), 2004, 470–507.
- [7] Brodsky, A., Sagiv, Y.: Inference of Monotonicity Constraints in Datalog Programs, *Principles of Database Systems*, ACM Press, 1989.
- [8] Bruynooghe, M., De Schreye, D., Krekels, B.: Compiling Control, *The Journal of Logic Programming*, **6**(1&2), 1989, 135–162.
- [9] Carlsson, M.: Freeze, Indexing and Other Implementation Issues in the WAM, *International Conference on Logic Programming* (J.-L. Lassez, Ed.), MIT Press, 1987.
- [10] Carlsson, M.: Personal communications with Mats Carlsson on the SICStus abstract machine, 2002–2004.
- [11] Cavedon, L.: Continuity, consistency, and completeness properties of logic programs, *International Conference on Logic Programming* (G. Levi, M. Martelli, Eds.), MIT Press, 1989.
- [12] Clark, K. L., McCabe, F. G., Gregory, S.: IC-Prolog Language Features, *Logic Programming* (K. L. Clark, S.-Å. Tärnlund, Eds.), Academic Press, 1982.
- [13] Dahl, V.: Two Solutions for the Negation Problem, *Workshop on Logic Programming* (S.-Å. Tärnlund, Ed.), 1980.
- [14] De Schreye, D., Decorte, S.: Termination of Logic Programs: The Never-Ending Story, *The Journal of Logic Programming*, **19&20**, 1994, 199–260.
- [15] De Schreye, D., Verschaetse, K.: Deriving Linear Size Relations for Logic Programs by Abstract Interpretation, *New Generation Computing*, **13**(2), 1995, 117–154.
- [16] Deransart, P., Ed-Dbali, A., Cervoni, L.: *Prolog: The Standard Reference Manual*, Springer-Verlag, 1996.
- [17] Genaim, S., King, A.: *Inferring Non-Suspension Conditions for Logic Programs with Dynamic Scheduling*, Technical Report 20-04, Computing Laboratory, 2004, See <http://www.cs.kent.ac.uk/pubs/2004/2008>.
- [18] Giacobazzi, R., Debray, S. K., Levi, G.: Generalized Semantics and Abstract Interpretation for Constraint Logic Programs, *The Journal of Logic Programming*, **3**(25), 1995, 191–248.
- [19] Hill, P. M., Lloyd, J. W.: *The Gödel Programming Language*, MIT Press, 1994.
- [20] Hoarau, S., Mesnard, F.: Inferring and Compiling Termination for Constraint Logic Programs, *Logic-based Program Synthesis and Transformation* (P. Flener, Ed.), 1559, Springer-Verlag, 1998.

- [21] Jaffar, J., Michaylov, S., Stuckey, P. J., Yap, R. H. C.: The CLP(R) Language and System, *ACM Transactions on Programming Languages and Systems*, **14**(3), 1992, 339–395.
- [22] Kowalski, R. A.: Algorithm = Logic + Control, *Communications of the ACM*, **22**(7), 1979, 424–436.
- [23] Lloyd, J.: *Foundations of Logic Programming*, Springer-Verlag, 1987.
- [24] Lüttringhaus-Kappel, S.: Control Generation for Logic Programs, *International Conference on Logic Programming* (D. S. Warren, Ed.), The MIT Press, 1993.
- [25] Marchiori, E.: Personal communication between Jonathan Martin and Elena Marchiori, 1996.
- [26] Marchiori, E., Teusink, F.: Termination of Logic Programs with Delay Declarations, *International Logic Programming Symposium* (J. W. Lloyd, Ed.), MIT Press, 1995.
- [27] Marchiori, E., Teusink, F.: Termination of Logic Programs with Delay Declarations, *The Journal of Logic Programming*, **39**(1–3), 1999, 95–124.
- [28] Martin, J. C., King, A.: Generating Efficient, Terminating Logic Programs, *Theory and Practice of Software Development* (M. Bidoit, M. Dauchet, Eds.), 1214, Springer-Verlag, 1997.
- [29] Martin, J. C., King, A.: On the Inference of Natural Level Mappings, *Program Development in Computational Logic* (M. Bruynooghe, K.-K. Lau, Eds.), 3049, Springer-Verlag, 2004.
- [30] Mesnard, F.: Towards Automatic Control for CLP( $\mathcal{X}$ ) Programs, *Logic Programming Synthesis and Transformation* (M. Proietti, Ed.), 1048, Springer-Verlag, 1995.
- [31] Mesnard, F., Ruggieri, S.: On Proving Left Termination of Constraint Logic Programs, *ACM Transactions on Computational Logic*, **4**(2), 2003, 207–259.
- [32] Naish, L.: *Negation and Control in Logic Programs*, Springer-Verlag, 1986.
- [33] Naish, L.: Coroutining and the Construction of Terminating Logic Programs, *Australian Computer Science Communications*, **15**(1), 1993, 181–190.
- [34] Pedreschi, D., Ruggieri, S.: Bounded Nondeterminism of Logic Programs, *International Conference on Logic Programming* (D. De Schreye, Ed.), MIT Press, 1999.
- [35] Pedreschi, D., Ruggieri, S., J.-G. Smaus: Classes of Terminating Logic Programs, *Theory and Practice of Logic Programming*, **2**(3), 2002, 369–418.
- [36] Rodríguez-Carbonell, E., Kapur, D.: An Abstract Interpretation Approach for Automatic Generation of Polynomial Invariants, *Static Analysis Symposium* (R. Giacobazzi, Ed.), 3148, Springer-Verlag, 2004.
- [37] SICS: *SICStus Prolog User's Manual*, 2005, See <http://www.sics.se/sicstus>.
- [38] Simon, A., King, A., Howe, J. M.: Two Variables per Linear Inequality as an Abstract Domain, *Logic Based Program Synthesis and Transformation* (M. Leuschel, Ed.), 2664, Springer-Verlag, 2003.
- [39] Takeuchi, A., Furukawa, K.: Bounded Buffer Communication in Concurrent Prolog, *New Generation Computing*, **3**(2), 1985, 145–155.
- [40] van Emden, M. H., de Lucena Filho, G. J.: Predicate Logic as a Language for Parallel Programming, *Logic Programming* (K. L. Clark, S.-Å. Tärnlund, Eds.), Academic Press, 1982.
- [41] Van Leeuwen, J., Ed.: *Handbook of Theoretical Computer Science: Volume B*, Elsevier, 1990.
- [42] Verschaetse, K., De Schreye, D., Bruynooghe, M.: Generation and Compilation of Efficient Computation Rules, *International Conference on Logic Programming* (D. H. D. Warren, P. Szeredi, Eds.), MIT Press, 1990.
- [43] Vieille, L.: Recursive Query Processing: The Power of Logic, *Theoretical Computer Science*, **69**(1), 1989, 1–53.